

# FROM FAB TO TOKEN

THE STATE OF THE AI INFRASTRUCTURE MARKET

---

QCON AI BOSTON · APRIL 2026

JORDAN NANOS · MEMBER OF TECHNICAL STAFF · SEMIANALYSIS

[boston.qcon.ai/presentation/boston2026/fab-token-state-market](https://boston.qcon.ai/presentation/boston2026/fab-token-state-market)

# INTRODUCTION

Who is SemiAnalysis?

SemiAnalysis is the best-of-breed Semiconductor and AI Research boutique. Our coverage spans the entire supply chain from semiconductor fabrication essentials to cutting-edge AI models, software, and the infrastructure that runs them.

## Product-first approach

Analysts have extensive operating and research experience in the semiconductor and AI industry. We talk directly to engineers, buyers, and operators to understand the technology itself.

## Connecting the dots

No other independent research house spends most of its time in the field seeking out public data points, applying technical expertise, and deriving actionable insights.

## Whole-stack coverage

From wafers, CoWoS and HBM to networking, CPO, datacenter capacity, hyperscaler and neocloud server TCO to tokenomics, we work as one global team across the entire AI supply chain.

**280,000+**

newsletter subscribers

**11**

integrated industry models

**6,000+**

data centers tracked globally

**100+**

conferences attended annually

# WHAT IS THE CHOKEPOINT?

This talk covers four things the market cares about. Critically, we assume models and talent are not the bottleneck, they keep improving.

01

## **Chips**

CoWoS, HBM, and wafer starts restrict how many accelerators reach customers each quarter.

02

## **Datacenters**

Multi-GW campuses with mandatory liquid cooling, 800VDC, and behind-the-meter gas turbines.

03

## **System Performance**

ClusterMAX for goodput and productivity, InferenceX for accelerators and their impact on margins.

04

## **End User Demand**

Tokenomics: ChatGPT MAU growth, Claude Code vs Codex, open-source impact, free vs paid monetization, and API revenue.

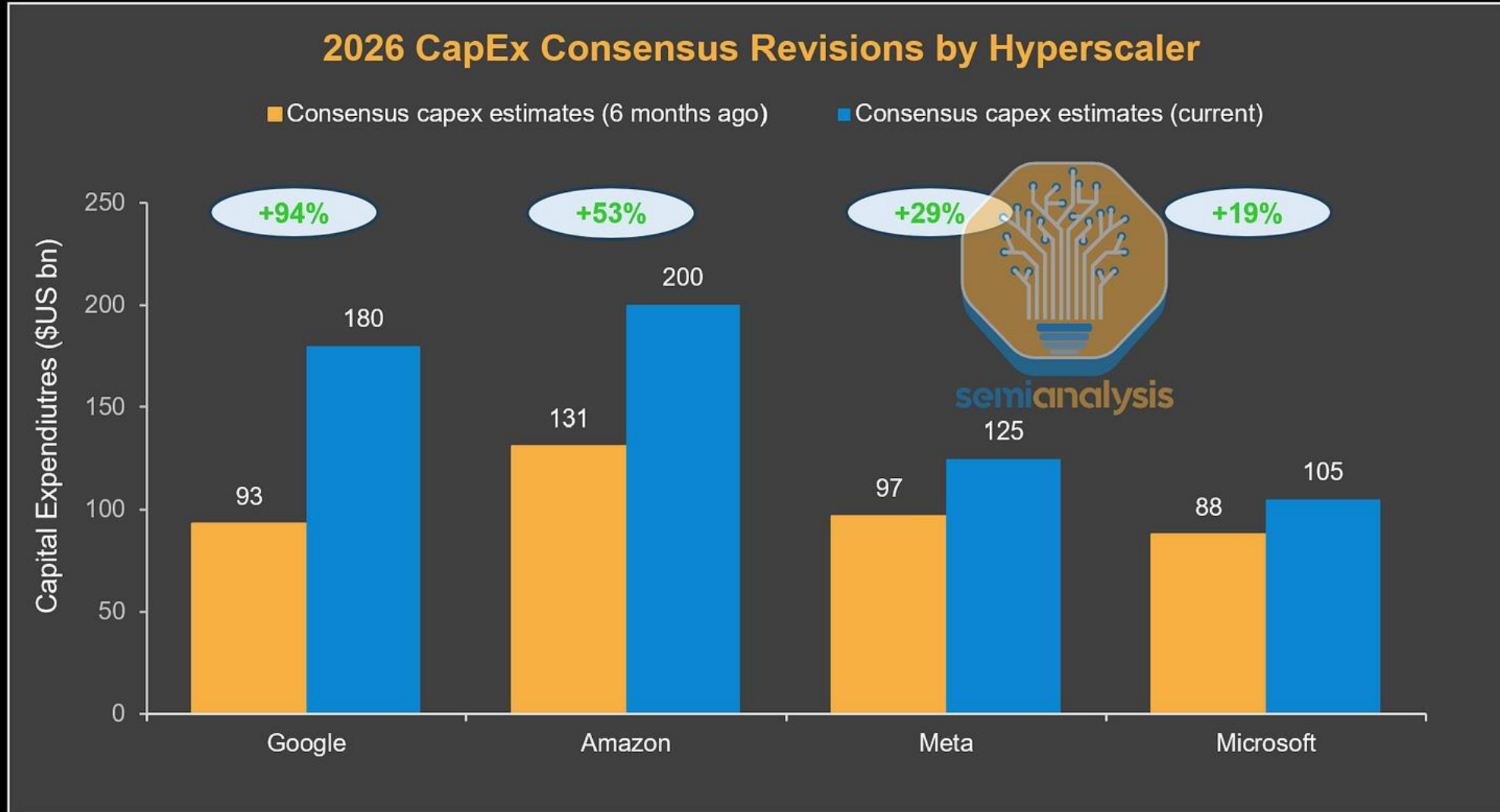
PART 01

# Chips

---

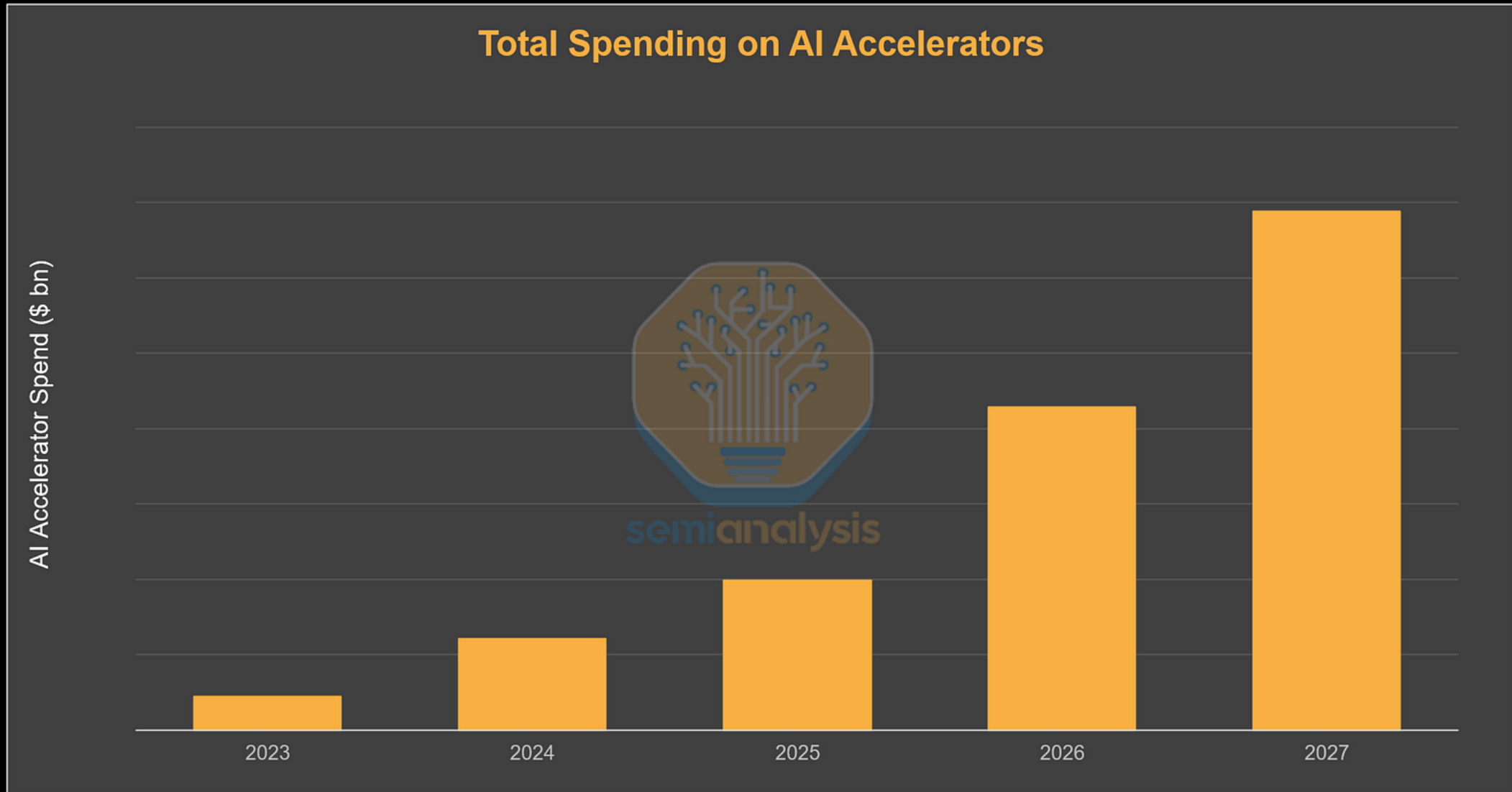
The accelerator market is gated by advanced packaging and HBM

# HYPERSCALERS ARE SPENDING A LOT OF MONEY ON CHIPS

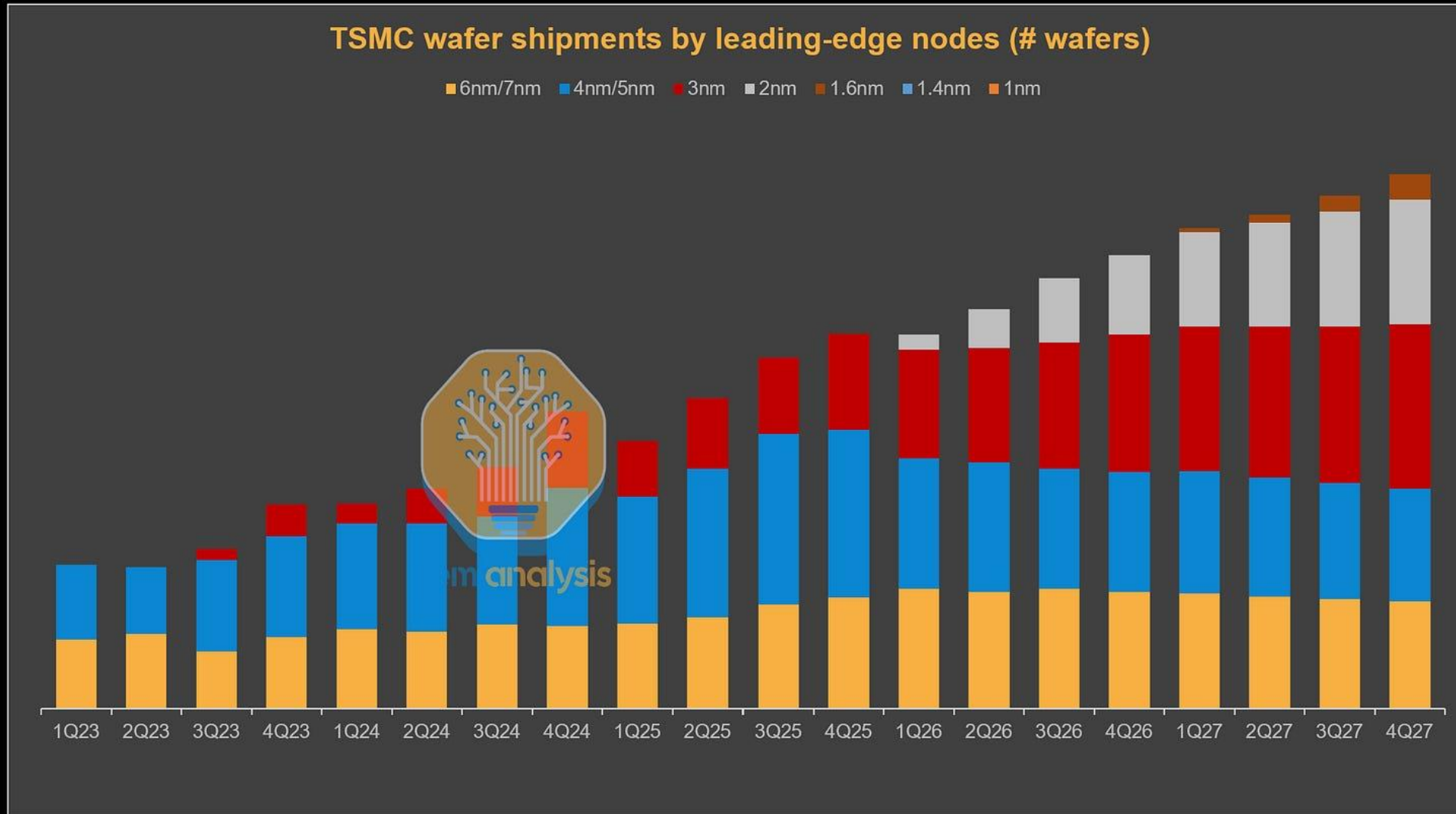


\*Microsoft capex does not include capital leases.

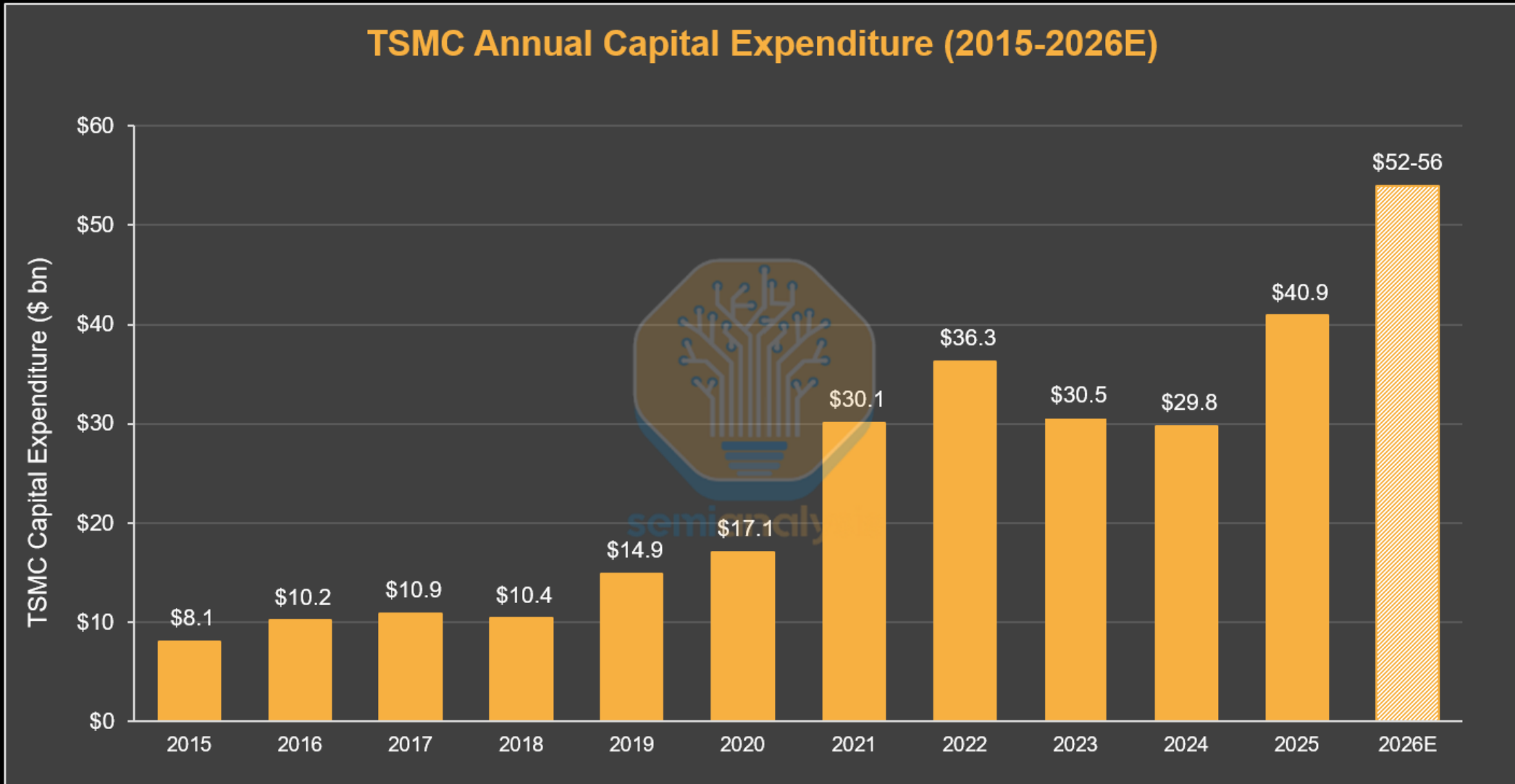
# HYPERSCALERS ARE SPENDING A LOT OF MONEY ON CHIPS



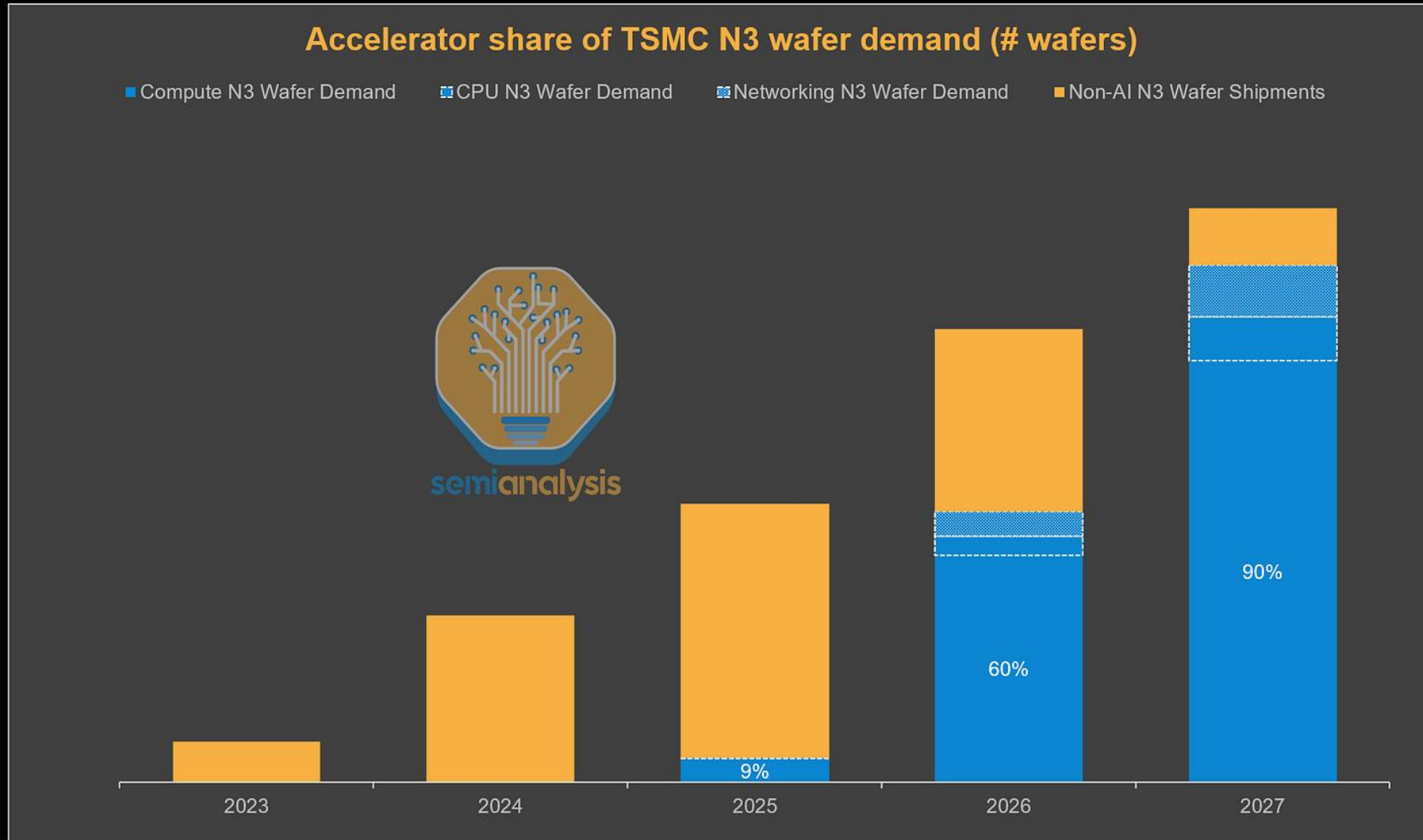
# TSMC IS NOT KEEPING UP WITH DEMAND GROWTH



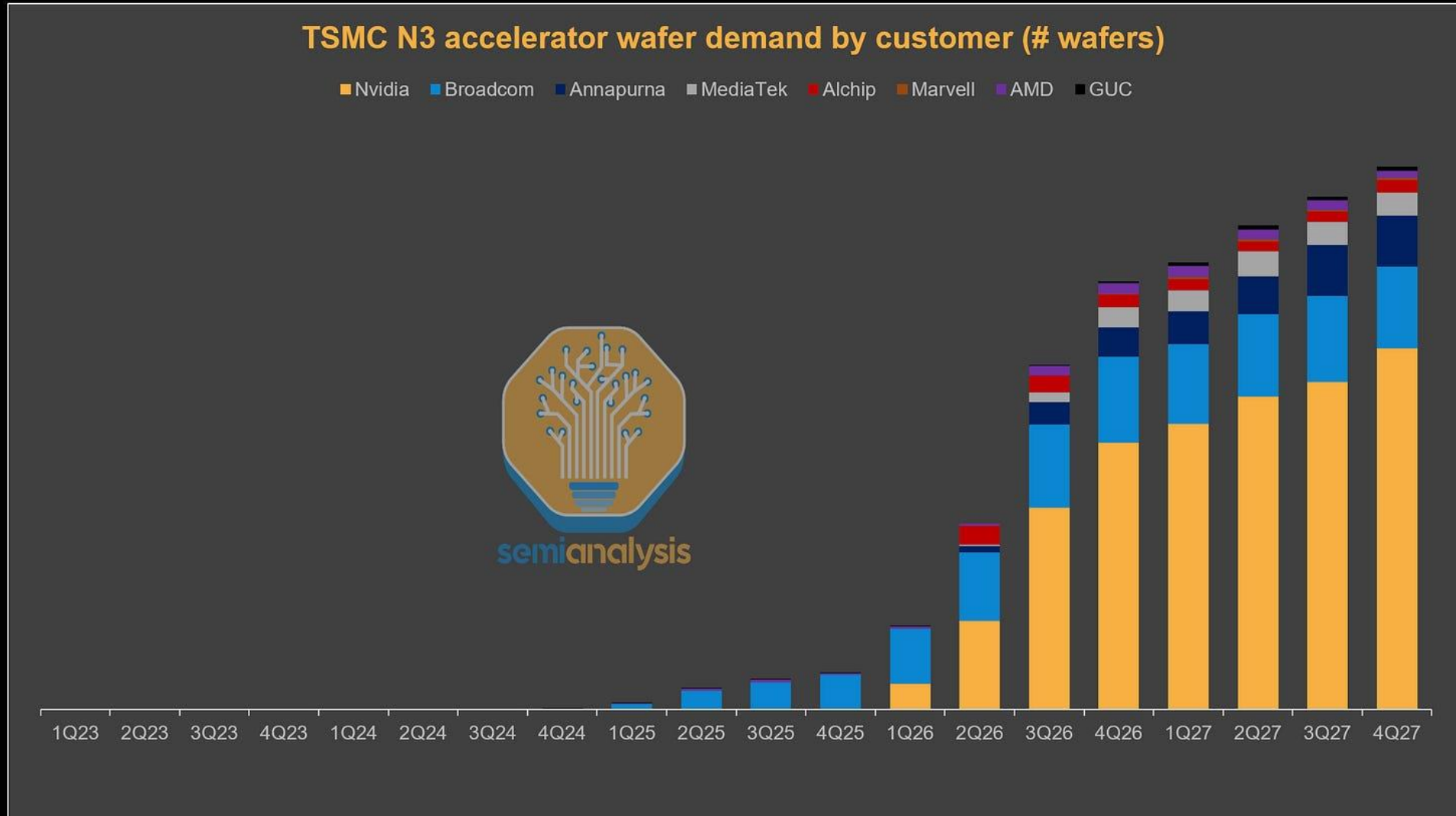
# TSMC IS NOT KEEPING UP WITH DEMAND GROWTH



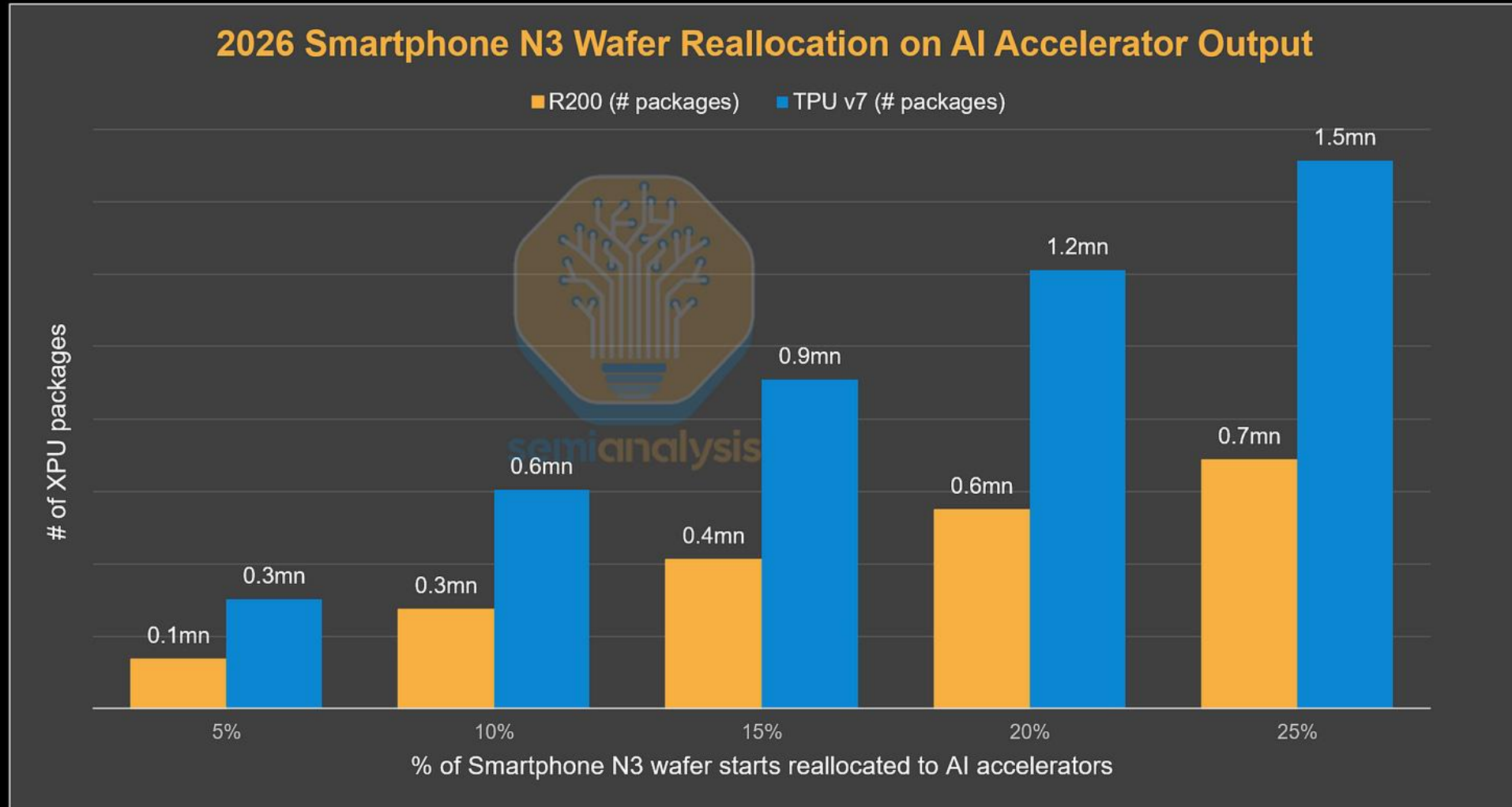
# TO BE CLEAR, ALL OF THE GROWTH IS BECAUSE OF AI



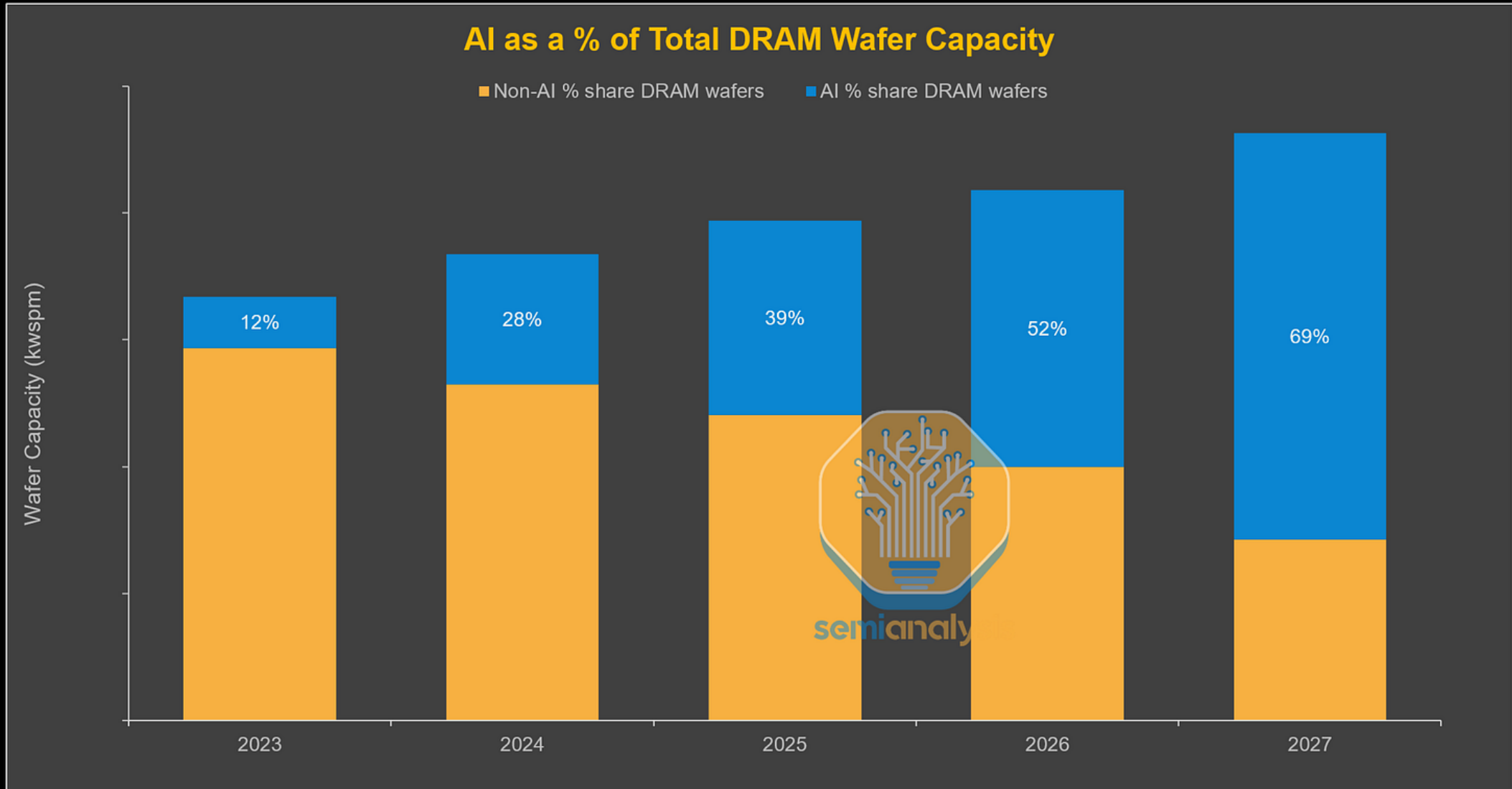
# TO BE CLEAR, ALL OF THE GROWTH IS BECAUSE OF AI



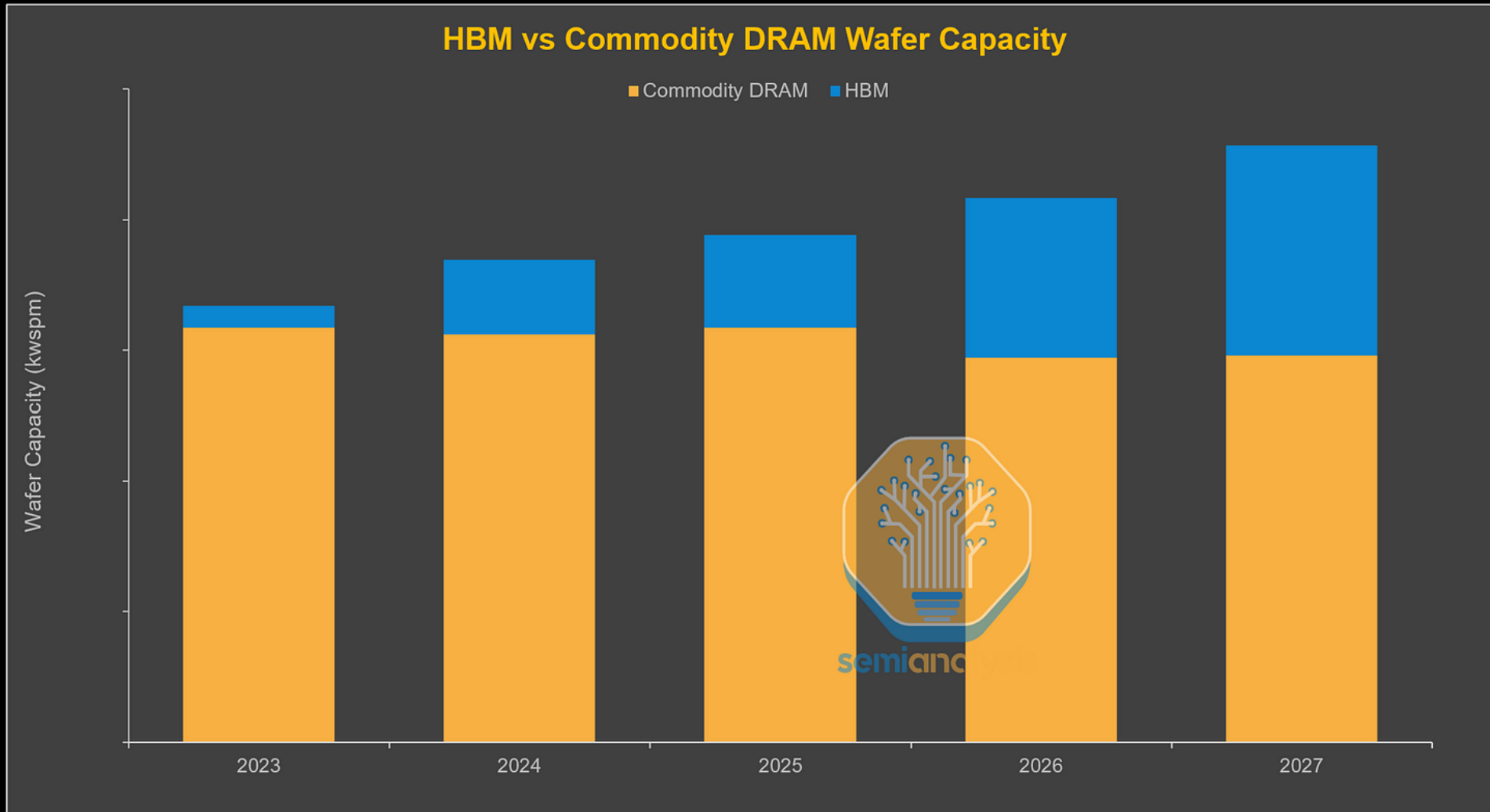
# AI GROWTH IS COMING AT THE EXPENSE OF SMARTPHONES



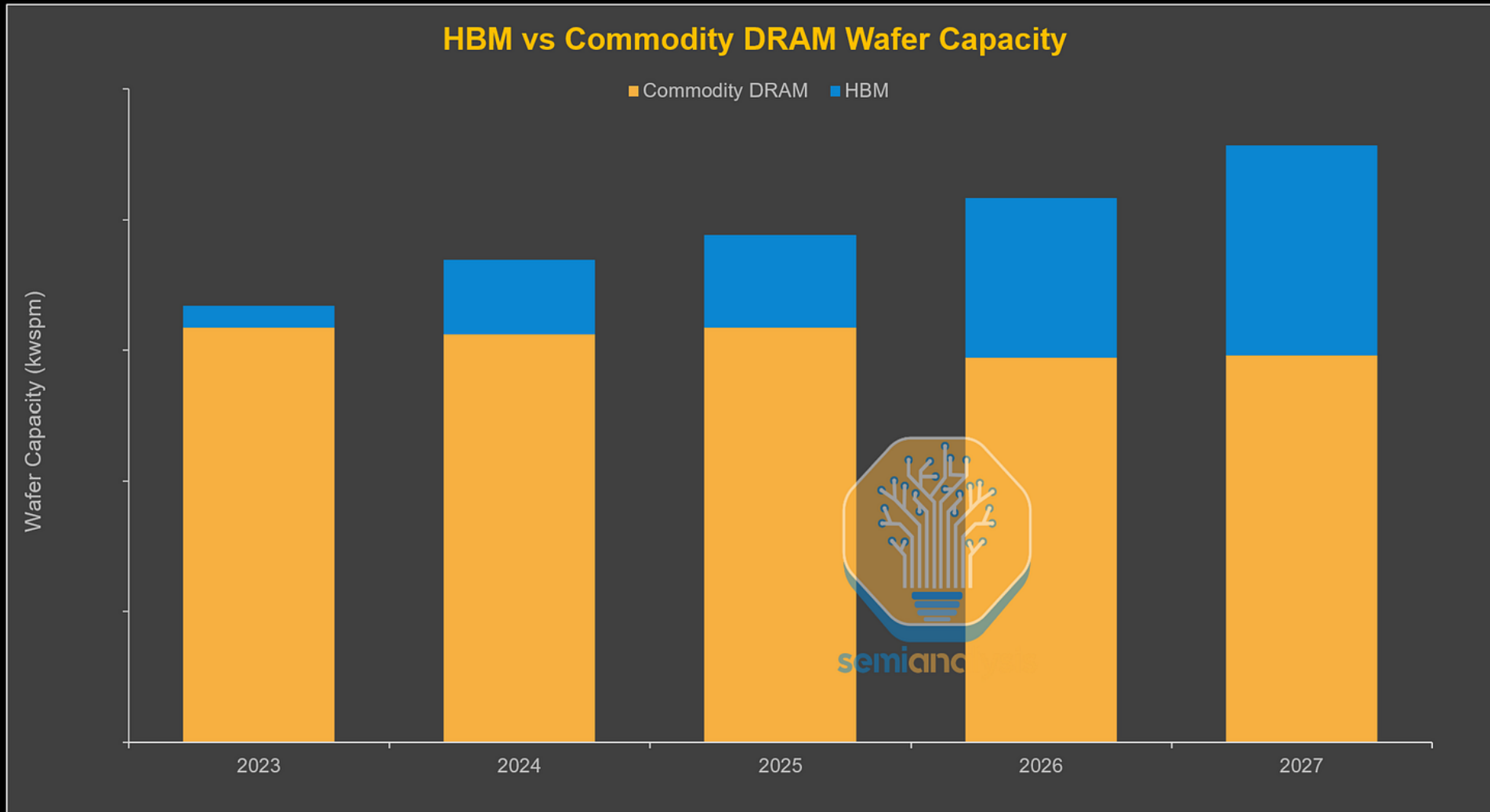
# AI IS ALSO COMING FOR MEMORY



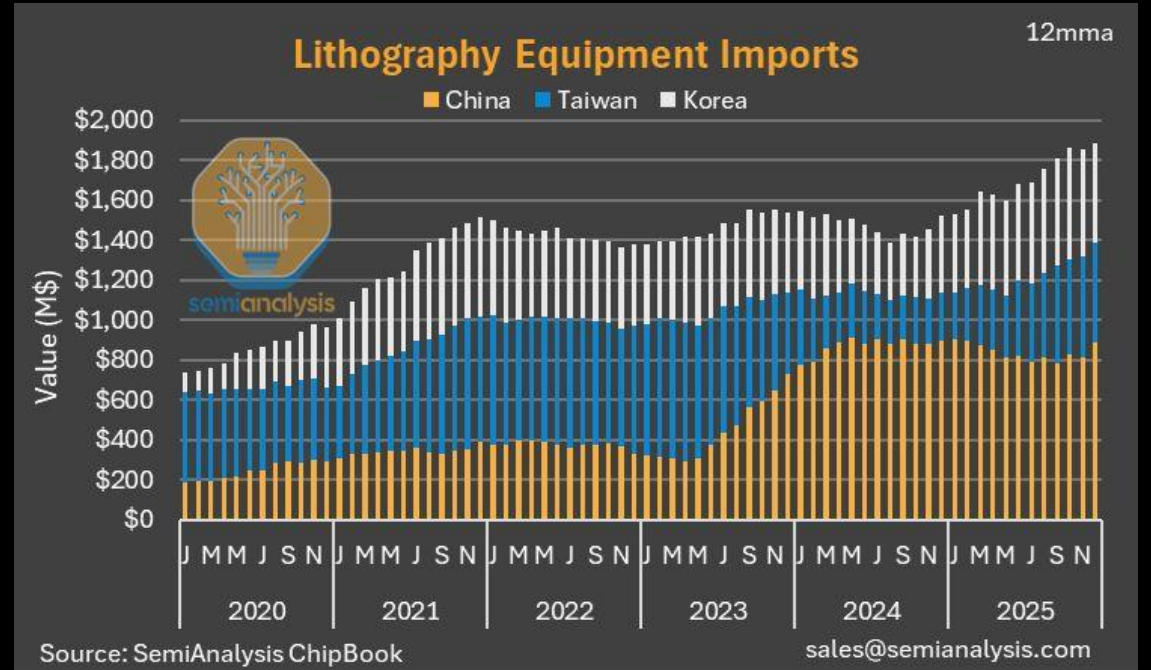
# AI IS ALSO COMING FOR MEMORY



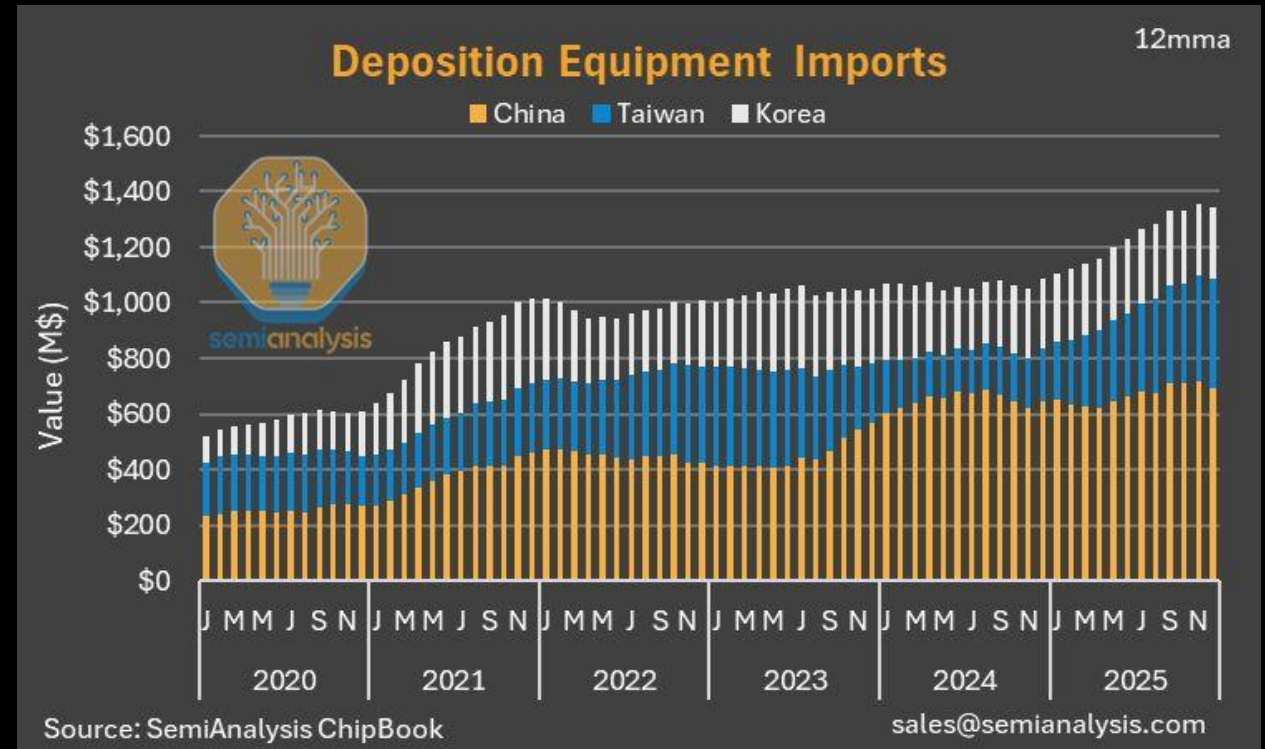
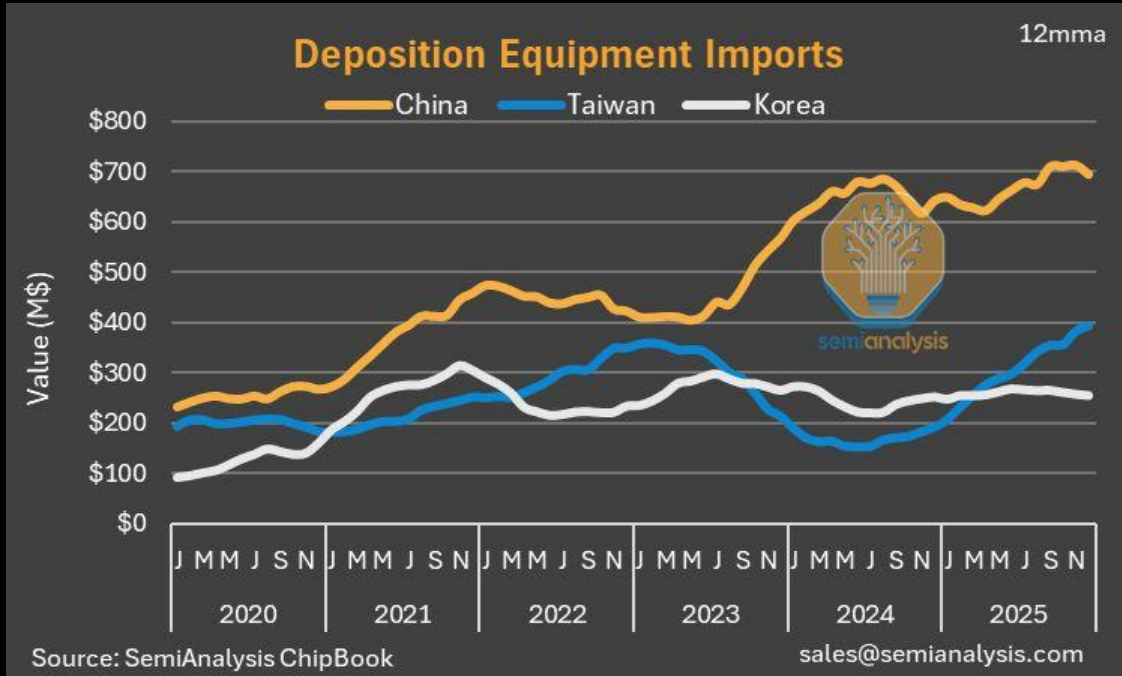
# AI IS ALSO COMING FOR MEMORY



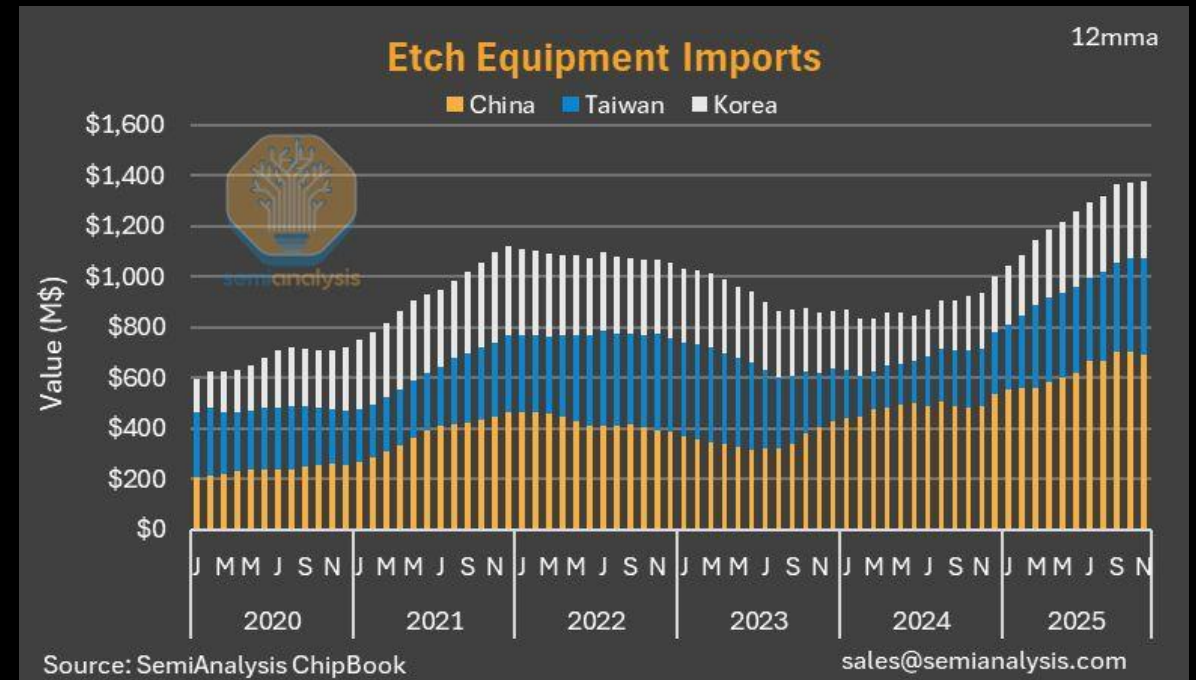
# Result: Wafer Fab Equipment (WFE) Spend Is Increasing



# Result: Wafer Fab Equipment (WFE) Spend Is Increasing



# Result: Wafer Fab Equipment (WFE) Spend Is Increasing



# NEW GPUS ARE COMING

NVIDIA Extreme Co-Design Delivering X-Factors Every Year  
From Chips to Racks to AI Factories



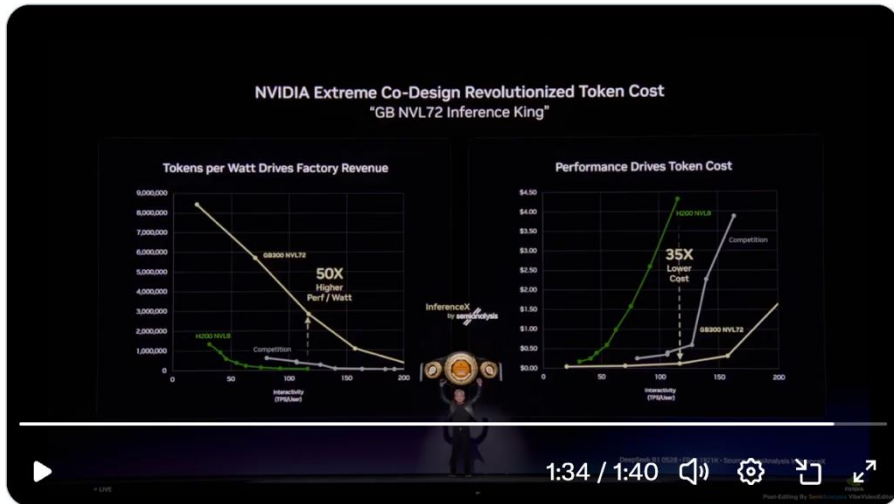
# SO FAR, PERFORMANCE IMPROVEMENTS HAVE BEEN IMMENSE

← Post



SemiAnalysis  
@SemiAnalysis\_

At GTC 2024, Jensen said that GB200 NVL72 was 35x faster than Hopper. Nobody believed it and thought it was classic fake Jensen Math. When we tested the performance of it, it wasn't just 35x faster, it was over 50x times faster even against a strong Hopper baseline with all of the inference optimization composed together like MTP, Disagg prefill, wideEP, etc. View the nuanced results at InferenceX dot com.



3:07 PM · Apr 18, 2026 · 60.8K Views

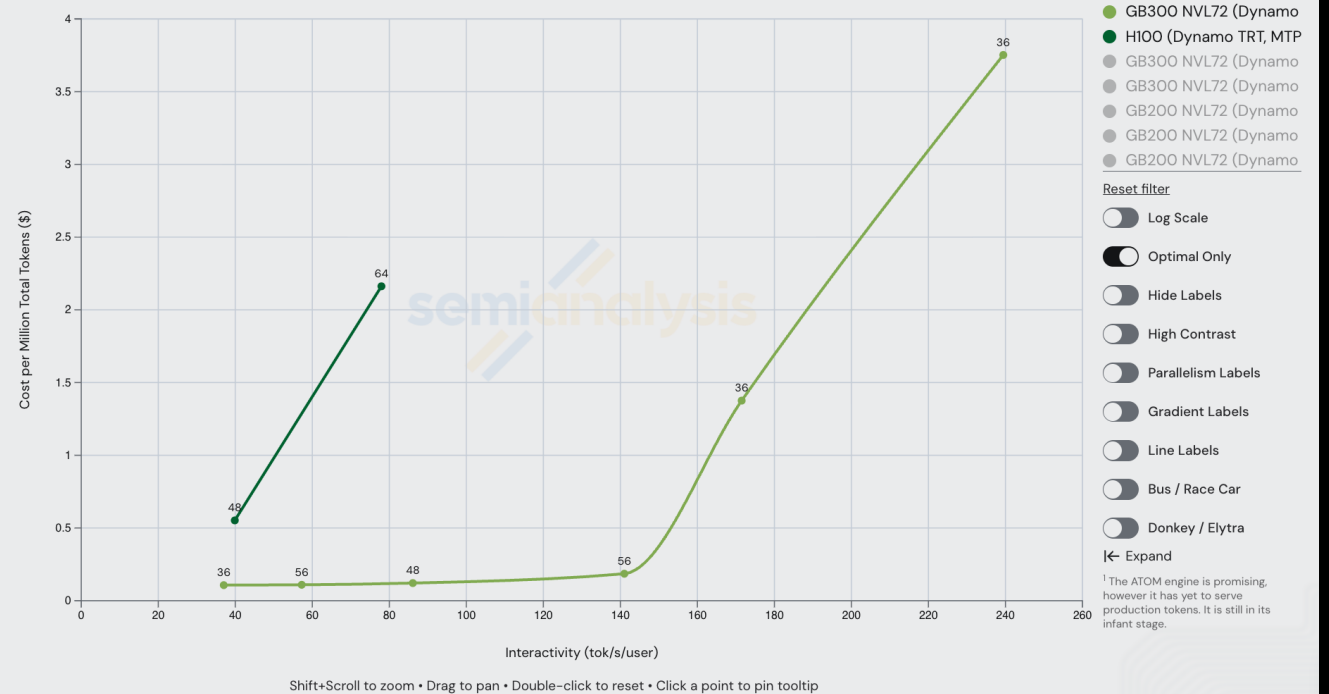
## Cost per Million Total Tokens (Owning – Neocloud Giant) vs. Interactivity

DeepSeek R1 0528 671B • FP8 • 8K / 1K • Source: SemiAnalysis InferenceX™ • Updated: 05/28/2026

TCO \$/GPU/hr: H100: 1.69 H200: 1.74 B200: 2.34 B300: 2.808 GB200: 2.75 GB300: 3.3 MI300X: 1.4 MI325X: 1.59 MI355X: 1.9

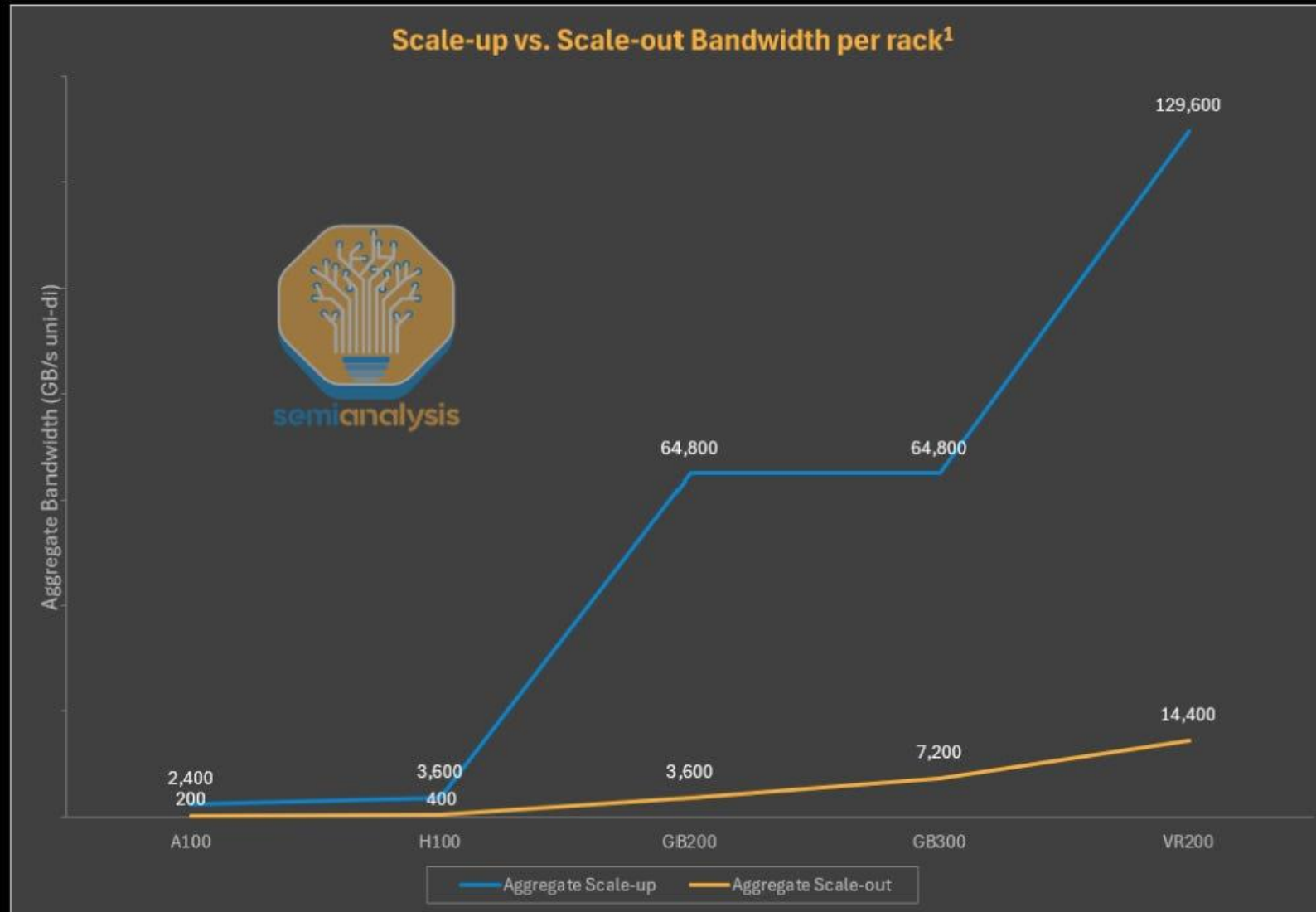
Source: [SemiAnalysis Market August 2025 Pricing Surveys & AI Cloud TCO Model](#)

**Note:** Disaggregated inference configurations (e.g., MoRI SGLang, Dynamo TRT) calculate cost per decode GPU or per prefill GPU, rather than per total GPU count. This makes direct cost comparison with aggregated configs not an apples-to-apples comparison.



\$2/Mtok on H100, \$0.10/Mtok on GB300 NVL72 (source: [inference.com](https://inference.com))

# NETWORKING IS THE PERFORMANCE BOTTLENECK



Notes: (1) Each rack has 8 GPUs for A100 and H100; each rack has 72 packages for GB200, GB300 and VR200

# CPO IS COMING



# COPPER TODAY, OPTICS TOMORROW

DAC → ACC → AEC → optical — each step trades reach, power, and cost

DAC	ACC	AEC	Optical / Kyber
<b>Passive copper</b>	<b>Active copper</b>	<b>Active electrical cable</b>	<b>Co-packaged + pluggable</b>
Reach <2 m reach	Reach 3–5 m	Reach 5–7 m	Reach 10–100+ m
Power 0 W per port	Power 1–2 W	Power 3–5 W	Power 8–15 W (LR/FR)
Where In-rack	Where Top-of-rack	Where Adjacent racks · NVL72 backplane	Where Scale-out · cross-row

NVL72 stays copper. NVL576 / Kyber-class designs need optical reach + co-packaged optics.

# CPO IS A POWER + RELIABILITY BET

Co-packaged optics cuts switch power and reach — but MTTF is the open question

## 22%

**less switch power**

CPO removes pluggable DSP + retimers — frees rack thermal budget for more compute.

## 20%

**lower \$/Gbps at scale**

Higher initial BOM, but lower lifetime energy + cable cost across 100k-port fabrics.

## MTTF

**is the question**

Field MTBF data is still thin; pluggables remain the safer choice for now.

### Vendor positioning

<b>NVIDIA</b>	Quantum-X Photonics + Spectrum-X Photonics — first 800G/1.6T CPO switches
<b>Broadcom</b>	Bailly / Tomahawk-6 CPO — pursued, then re-scoped as pluggable-first
<b>Marvell + partners</b>	CPO PHY IP licensing into custom silicon
<b>Status</b>	Pluggable optics still ship the volume — CPO penetrates only at the highest radix tiers

PART 02

# Datacenters

---

Megawatts per rack, gigawatts per campus — the grid is the new bottleneck.

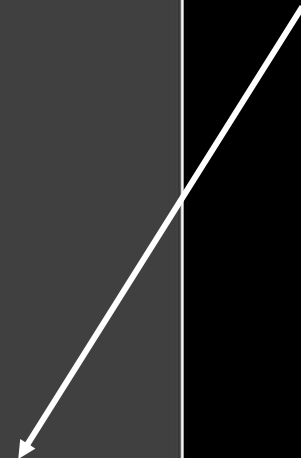
# WHERE DO ALL THESE CHIPS GO?

## Global DC Critical IT Demand vs Supply (GW) - excluding China

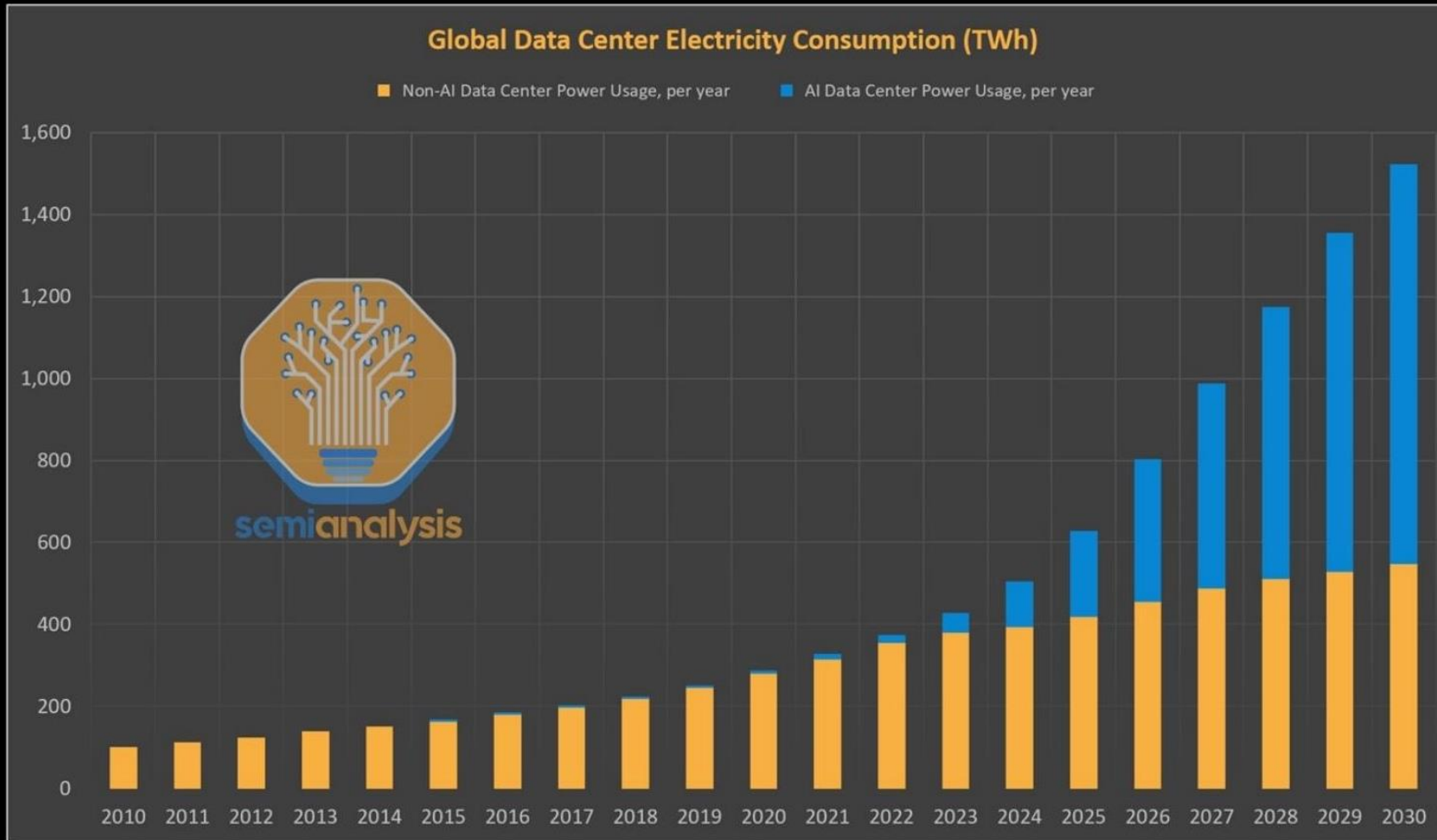
■ AI + Non-AI Data Center Demand   ■ Total Self-Build + Colocation Capacity   ■ Total Surplus / (Deficit)



Expanding datacenters and power is much easier than building new fabs.

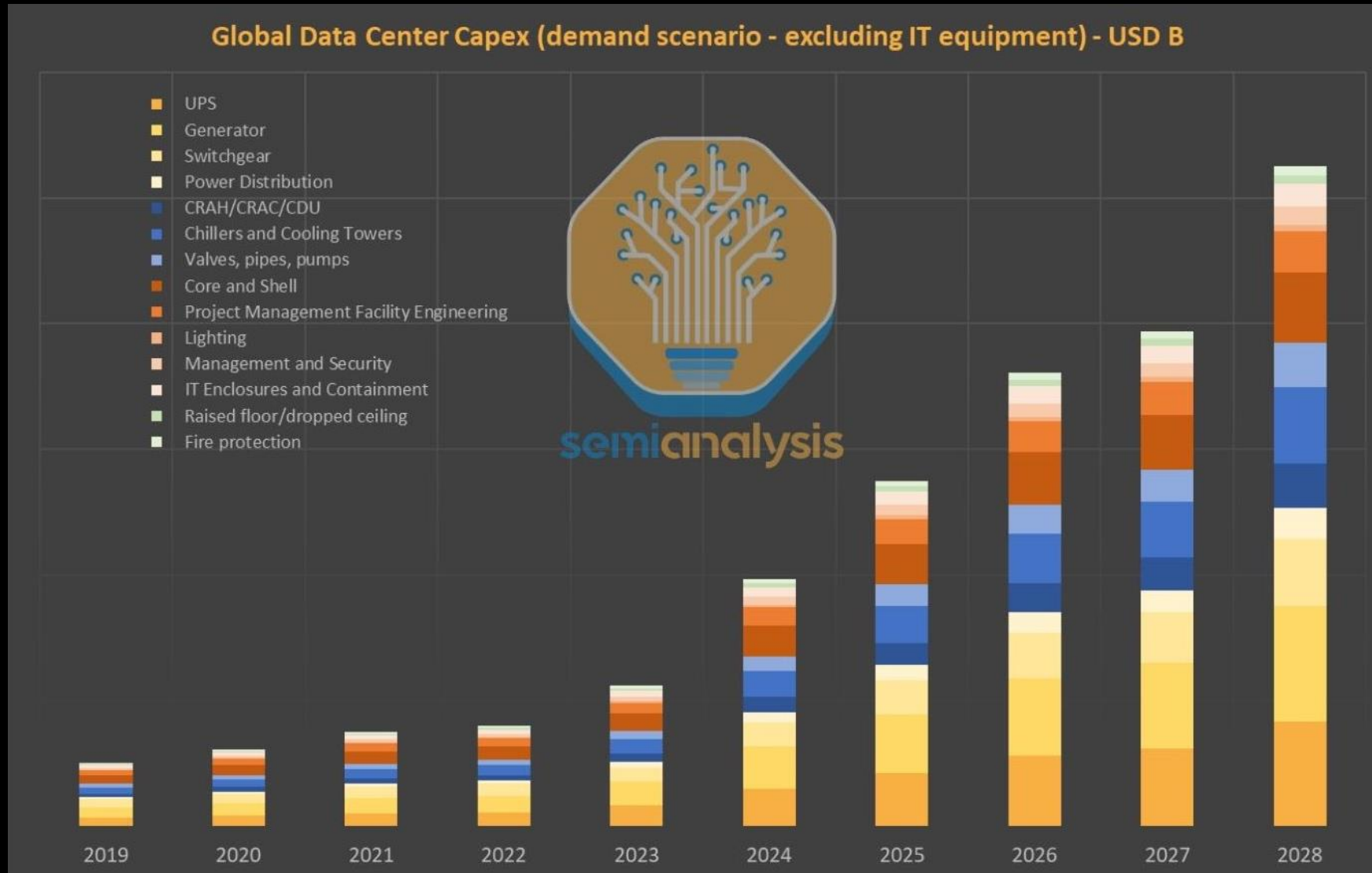


# Industry-wide load growth – GW additions by region



Forecast critical IT load additions by region and customer (hyperscaler, neocloud, enterprise) through 2030.

# Industry-wide load growth – GW additions by region



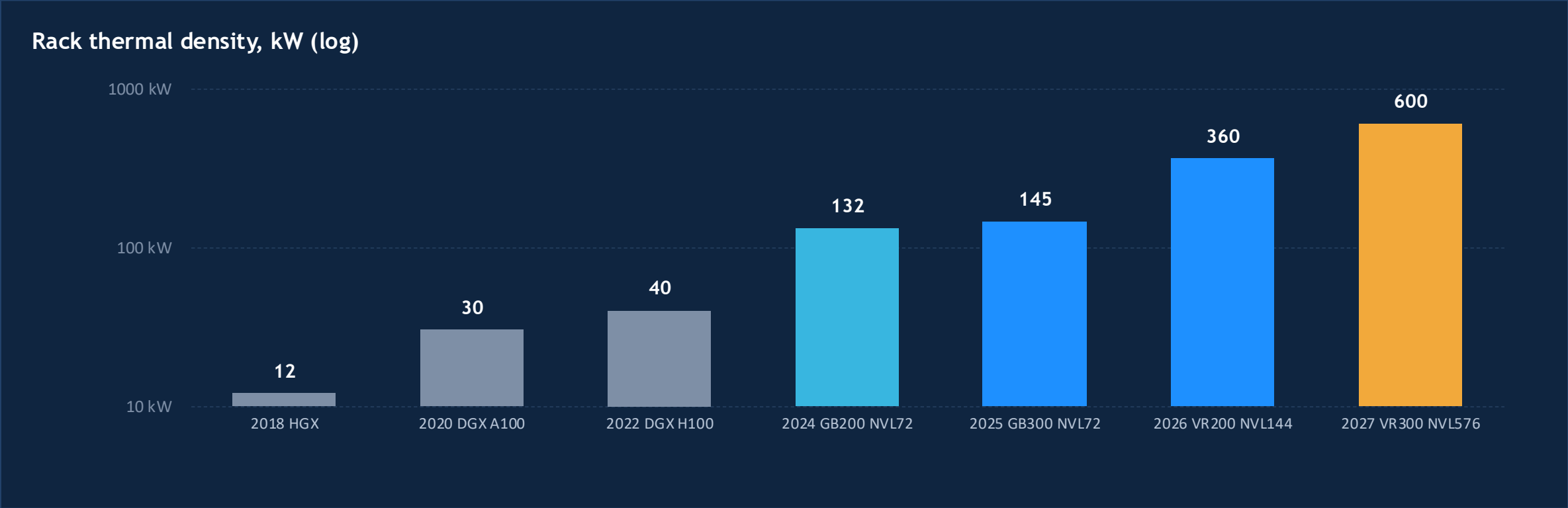
# YOU CAN GO SEE IT FOR YOURSELF



(pictures from Friday last week)

# RACK POWER HAS GONE EXPONENTIAL

From 12 kW air-cooled to 600+ kW liquid-cooled in one product generation



50× density growth in 9 years – air cooling cannot follow this curve.

Source: [NVIDIA datasheets](#), [SemiAnalysis Datacenter Model](#)

# THESE GPUS ARE COMING

← **Michael Dell**  
6,034 posts



**Michael Dell** 🇺🇸

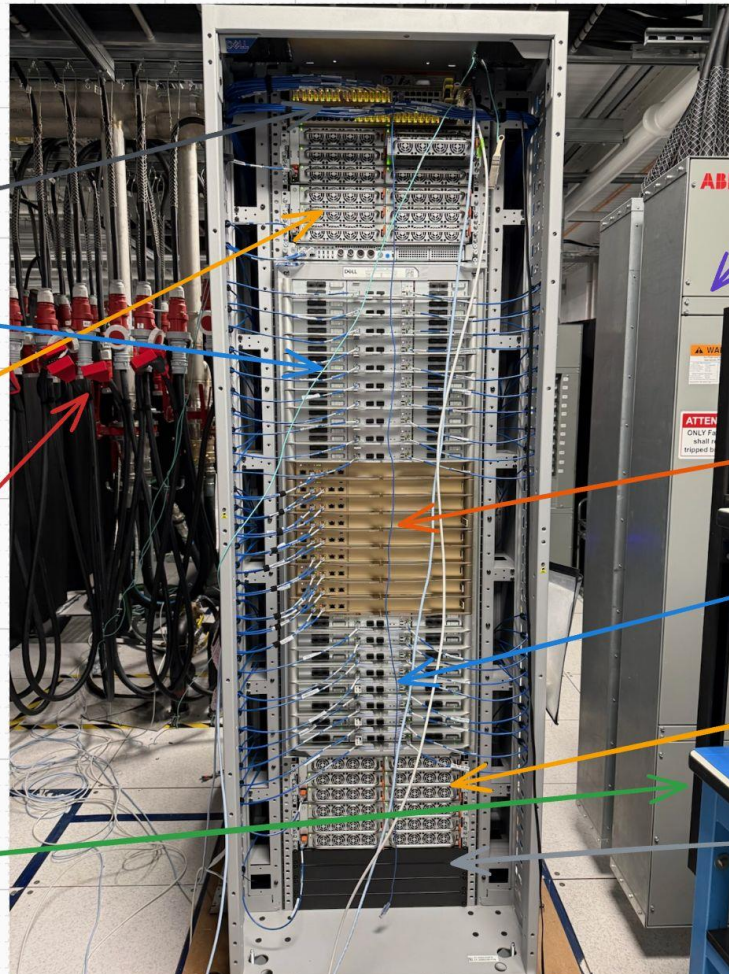
@MichaelDell

Husband, Father, Grandfather  
Technologies Grateful 🇺🇸

📅 Entrepreneur 📍 Australia  
📅 Joined July 2009 >

4,166 Following 2.9M Followers

NVIDIA Vera Rubin VR200 NVL72 Oberon Rack



Management switches

9 compute trays (top half)  
4x R200 GPU pkgs per tray

Power shelves (top half) 12x ORVS 33 kW IU  
6x 5.5 kW PSUs, 48-54 VDC, 660-700A

Power whips

Crash cart

ABB Remote Power Panel / LV switchgear

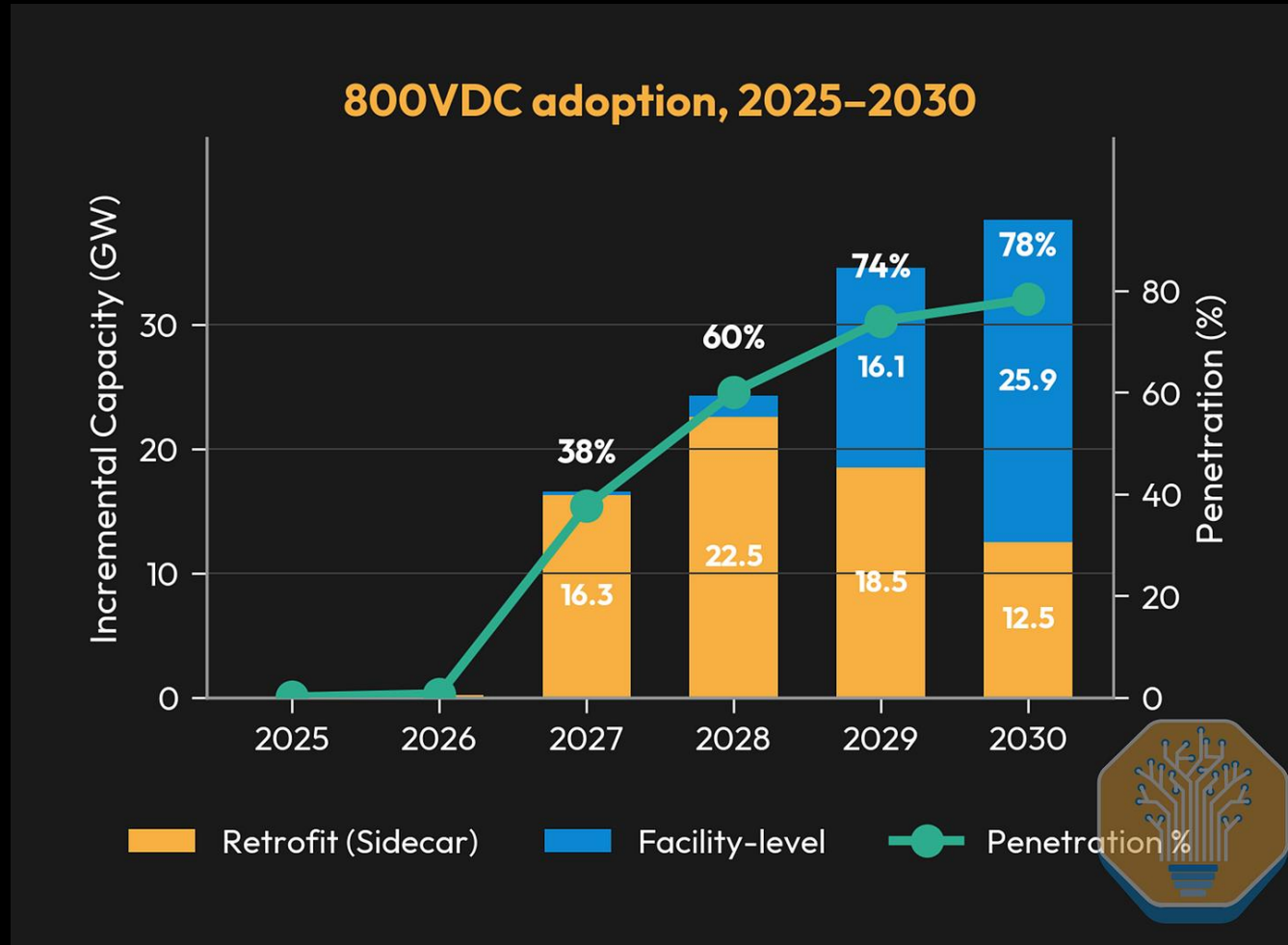
NVSwitch trays spray painted gold

9 compute trays (bottom half)  
18 trays total x 4 R200 = 72 packages

Power shelves (bottom half)  
440kW total capacity (N+N)

Blanks

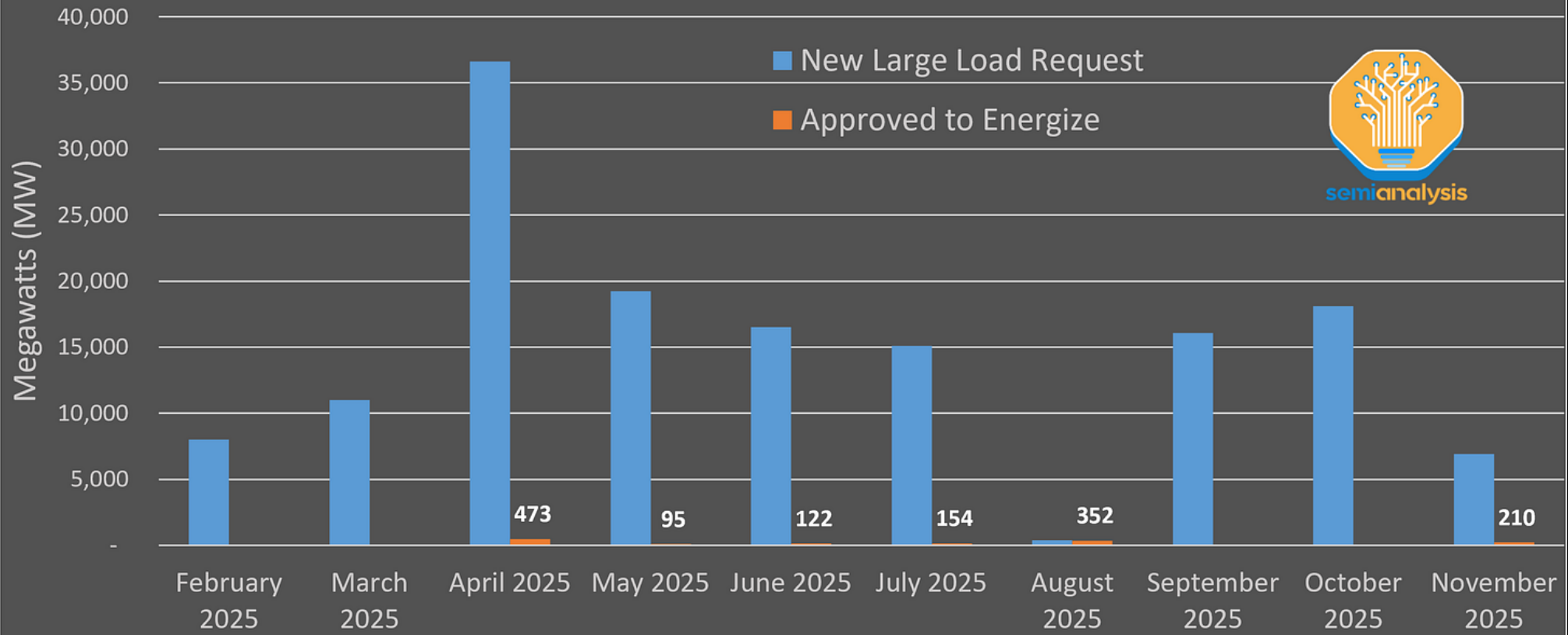
# THE RESULT IS 800VDC



# BUT THE GRID IS SOLD OUT

## Texas: Large Loads Requests vs Approved To Energize (MW)

The System is Sold Out



# THE SOLUTION: BEHIND-THE-METER GAS



xAI Colossus 1 - twelve SMT-130 turbines (198MW)  
Source: SemiAnalysis Datacenter Industry Model



xAI Colossus 2 - seven Titan 350 turbines (266 MW)  
Source: SemiAnalysis Datacenter Industry Model

# Construction tracking – satellite + permits



October 24, 2024

March 21, 2025



March 3, 2024

December 23, 2024



April 21, 2024

April 8, 2025



October 24, 2024

March 25, 2025



*Satellite imagery + permit + interconnect data gives us a real-time view of construction progress across every major AI campus.*

# Are AI Datacenters Increasing Electric Bills for American Households?

*PJM capacity surge narrows the bill gap: ERCOT-PJM spread compresses from 18/mo(2025)to5/mo (2026)*

*Source: EIA, Monitoring Analytics SOM, Potomac Economics SOM, PJM BRA Reports, SemiAnalysis Estimates*

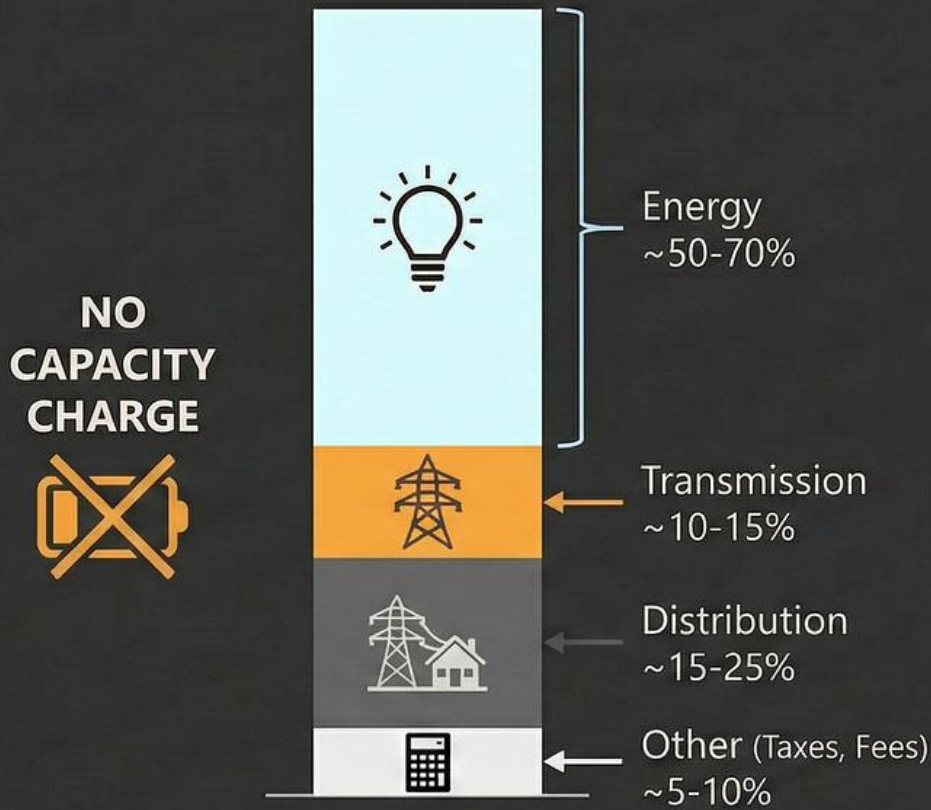
Annual Average Monthly Household Electric Bill Decomposition (\$/month)



# Are AI Datacenters Increasing Electric Bills for American Households?

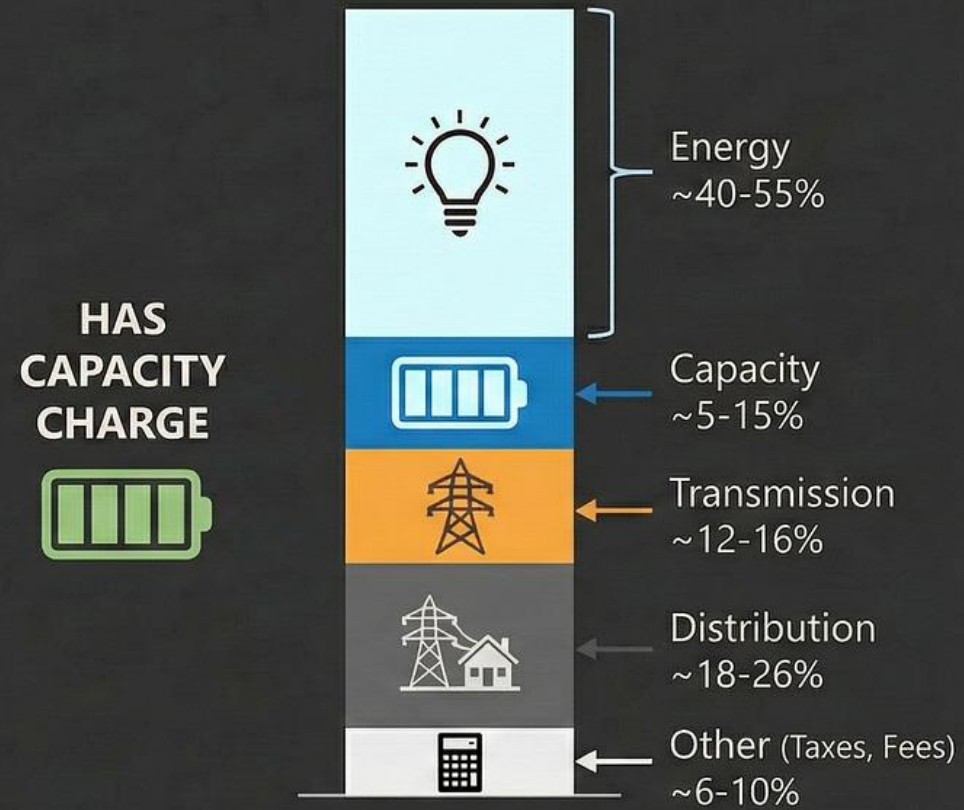
## Consumer Electricity Bill Breakdown

Bill Breakdown (Approx. Ranges)



Ranges are illustrative and can vary significantly based on provider, location, and time.

Bill Breakdown (Approx. Ranges)

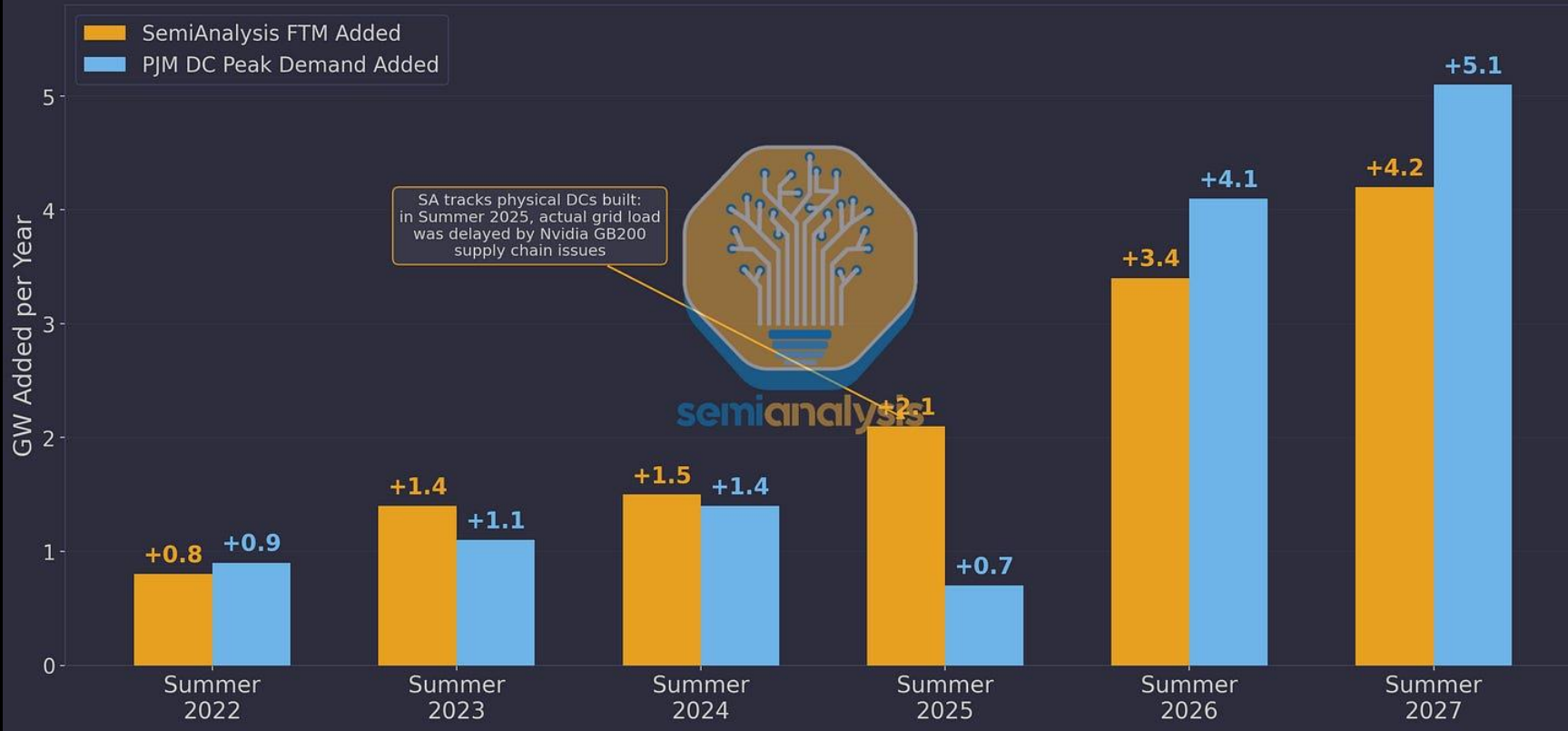


Ranges are illustrative and can vary significantly based on provider, location, and time.

# Energy Model – generation mix, interconnect, behind-the-meter

Source: SemiAnalysis Datacenter Model, Monitoring Analytics BRA Reports, PJM 2026 Load Forecast

## Annual DC Front-of-Meter Additions in PJM (GW/year)

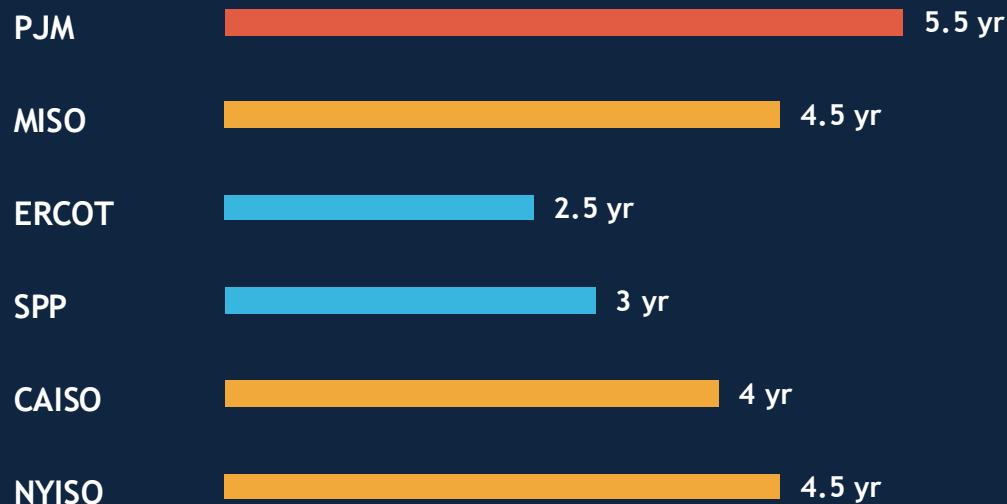


Power generation, transmission, and behind-the-meter capacity additions feeding AI data center demand through 2030.

# BEHIND-THE-METER IS NOW THE PATH

Grid interconnect queues stretch 4–7 years — speed-to-energy decides who trains first

## Typical interconnect-queue wait (years)



## Who is going behind-the-meter

<b>xAI Memphis</b>	Gas turbines on-site, ~150 MW added in months
<b>OpenAI Stargate</b>	Multi-GW build with co-located generation
<b>Meta Hyperion</b>	Louisiana, gas + storage hybrid
<b>Crusoe + OAI TX</b>	Pad-side gas — fast permit path
<b>Anthropic + AWS</b>	Nuclear PPAs for baseload, multi-year horizon

PART 03

# System Performance

---

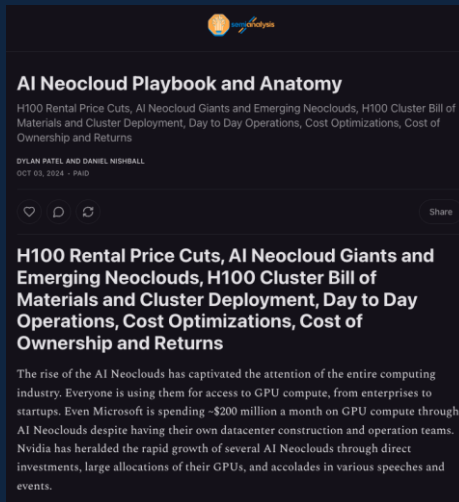
Two independent, public projects: ClusterMAX compares neoclods, InferenceX compares chips

# WHAT IS A NEOCLOUD?

The collage features four articles:

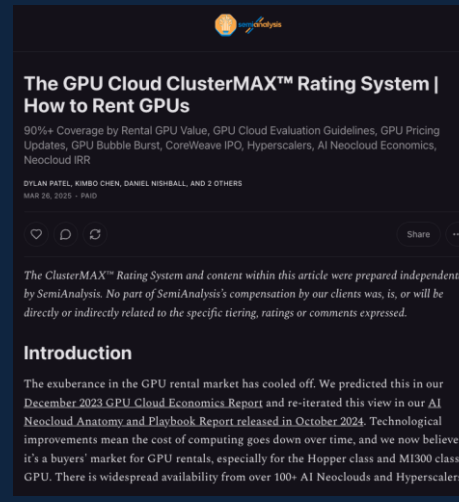
- Forbes:** "Neoclouds' Rise Reflects How AI Is Transforming The Cloud Market" by Matt Kimball, published March 25, 2026.
- EQUINIX Interconnections:** "THE INFRASTRUCTURE BEHIND AI: What Is a Neocloud?"
- McKinsey & Company:** "The evolution of neoclouds and their next moves" by Massimo Mazza, Pankaj Sachdeva, Suren Arutyunyan, and Tarik Alatovic, published November 19, 2025.
- Gartner:** "Gartner: Why neoclouds are the future of GPU-as-a-Service" by Mike Dorosh, published February 20, 2026. The article states that neocloud providers will capture around 20% of the \$267bn AI cloud market by 2030.

# HOW WE GOT HERE



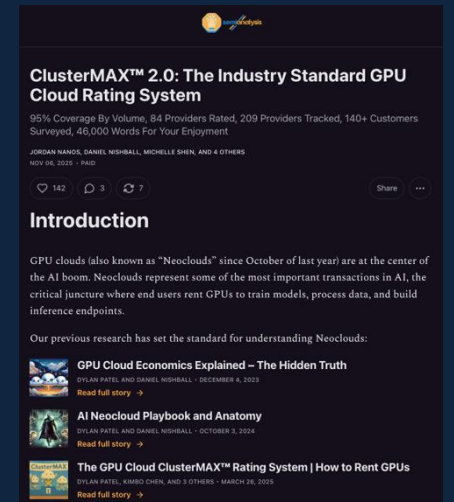
SemiAnalysis coins the term  
"Neocloud"

October 3, 2024



The first ClusterMAX rating system  
is released

March 26, 2025



ClusterMAX 2.0 is released

November 6, 2025



# WHAT MAKES A GOOD NEOCLOUD?

## Expectations

Explore our detailed expectations for different deployment models and infrastructure components

SLURM

Kubernetes

Standalone

Monitoring

Health Checks

## Criteria

Explore the 10 critical dimensions we use to evaluate GPU cloud providers

Security

Lifecycle

Orchestration

Storage

Networking

Reliability

Monitoring

Pricing

Partnerships

Availability

# EXAMPLE TESTS WE RUN

## Cluster audit, review of monitoring dashboard

### Common issues for [slurm](#):

- passwordless ssh to worker nodes, sudo access
- nvidia container toolkit installed and up to date
- ncu installed with hardware counters accessible
- high bw NICs/HCAs are named correctly and not conflicting, allowing NCCL to pick them up: mlx5\_0, mlx5\_1, etc.
- cuda, nvcc, hpc-x, nccl pre-installed from [hpc sdk](#)
- lmod installed and configured
- pyxis, enroot, docker installed and up to date
- topology.conf properly configured
- dcgm's [background health checks](#) enabled and plugged into slurm's [HealthCheckProgram](#)
- prolog and epilog scripts are lightweight, takes less than 30 seconds to get a node
- dashboard includes slurm job accounting data via [sacct](#) for resource usage, summaries by user and group

### Common issues for [kubernetes](#):

- passwordless ssh to worker nodes, sudo access
- nvidia container toolkit installed and up to date via GPUOperator or equivalent
- ncu installed with hardware counters accessible
- high bw NICs/HCAs are named correctly and not conflicting, allowing NetworkOperator to work with rdma: 1 pod spec (or equivalent)
- kubeconfig is a simple download
- helm access available without custom external authentication
- csi provider has ReadWriteMany volume support
- default StorageClass configured and functional, PVCs provision without hanging on helm install

# EXAMPLE TESTS WE RUN (CONT.)

## Real-world workloads

- torchtitan pretraining example: <https://github.com/pytorch/torchtitan>
- prime-rl example: <https://github.com/PrimeIntellect-ai/prime-rl>
- inferencex example: <https://github.com/SemiAnalysisAI/inferencex>
  
- Networking
  - nccl-tests: <https://github.com/NVIDIA/nccl-tests> (or [rccl-tests](#)) allreduce, allgather, alltoall, etc.
  - ib\_write\_bw and ib\_read\_bw from perftest: <https://github.com/linux-rdma/perftest>
  - stas00 [all\\_reduce\\_bench.py](#) (similar to nccl-tests, but launched from torch without mpi)
  
- Storage
  - fio: <https://github.com/axboe/fio> sequential read/write and random read/write IOPS
  - custom torch.save and dcp.save test
  - custom synthetic data generation test
  - time to run “import torch” on home directory

# EXAMPLE TESTS WE RUN (CONT.)

## Compute

- nvbandwidth: <https://github.com/NVIDIA/nvbandwidth> h2d, d2h, p2p
- stas00 [mamf-finder.py](#) custom GEMM benchmark from torch, triton, deepgemm

## Lifecycle

- synthetic download/upload speed (speedtest)
- hf download/upload (model files)
- container pull/push (ngc, ghcr, dockerhub)
- (uv) pip install torch
- add new users

## Reliability

- 8hr burn-in checking for thermal, power, ECCs, link flaps, storage mounts, etc.
- Xid error detection
- sudo reboot

# CLUSTERMAX RATINGS GO BEYOND HANDS ON TESTING

To accurately assess providers on reliability, performance, support and pricing at scale we interview hundreds of users

We also introduced a methodology for measuring the impact of downtime, aka “goodput” in our latest article:



## How Much Do GPU Clusters Really Cost?

Calculating Cluster TCO, The Real Impact of Downtime, The Grand Unifying Theory Of Goodput, and a ClusterMAX 2.1 Update

JORDAN NANOS, BRYAN SHAN, CHEANG KANG WEN, AND 2 OTHERS

APR 20, 2026 · PAID

123

2

9

Share

...

## Introduction: Rethinking the Total Cost of a GPU Cluster

Modern GPUs are unbelievably expensive. A single Blackwell GPU costs more than the average car, and uses more energy than a single family home. It is now common for unicorn startups to have thousands of these GPUs working for them, day and night. Many foundation model companies now spend an order of magnitude more money on GPUs than they do on employees. We know multiple companies spending over 80% of their initial funding on GPUs. Startup founders now have four important categories of spending to consider when building a financial plan for their company:

1. GPU clusters
2. Tokens
3. Employees
4. Everything else

# CLUSTERMAX RATINGS GO BEYOND HANDS ON TESTING

To accurately assess providers on reliability, performance, support and pricing at scale we interview hundreds of users

We also introduced a methodology for measuring the impact of downtime, aka “goodput” in our latest article:

Item	Qty	Gold-tier	Hyperscaler	Silver-tier
<b>GPU</b> \$/GPU-hr + GPU + Premium				
GB300 NVL72	5184	\$4	\$4	\$4
Orchestration		included	0 % incl	included
<b>Storage</b> \$/GiB-mo +				
Hot (Lustre/Weka)	500 TIB	\$0.035	\$0.0725	\$0.055
Cold (S3/Object)	10 PIB	\$0.01	\$0.02	\$0.015
<b>Network</b> \$/mo				
		\$0	\$19	\$0
<b>CPU</b> \$/vm-hr				
		\$0	\$3.318	\$0
<b>Support</b> % uplift				
		\$0	\$455,403	\$0
<b>Goodput</b> % uplift +				
		\$917,088	\$1,571,424	\$3,121,349
Goodput	ongoing	6.14 % incl	10.53 % incl	20.91 % incl
<b>Setup</b> one-time				
		\$0	\$14,996,587	\$33,333
<b>Debugging</b> \$/mo				
		\$0	\$8,333	\$128,583
<b>Subtotals</b> 3 years <input type="checkbox"/> Incl. goodput/setup/debug				
Monthly (amortized)		\$15,969,785/mo	\$17,631,823/mo	\$18,366,224/mo
36-month Total		\$574,912,276	\$634,745,636	\$661,184,050
Relative to Gold-tier		1.00x	1.10x	1.15x

Parameter	Gold-tier	Hyperscaler	Silver-tier
<b>SHARED</b>			
Cluster size (GPUs)	5184	5,184	5,184
Avg job size (j_size) (GPUs)	4096	4,096	4,096
Blast radius (b_radius) (GPUs)	64	64	64
<b>PER PROVIDER</b>			
GPU MTBF (GPU-hrs)	25000	25000	15000
Checkpoint freq (t_chkpt) (mins)	60	60	60
Failover time (t_failover) (mins)	5	5	5
Idle spare GPUs (GPUs)	32	0	0
<b>RESILIENCY</b>			
Repair/Replace	Fault Tolerant	Fault Tolerant	Checkpoint Restart
	TorchPass	HyperPod Ckptless	Hot spare (idle)
Time to identify failure (mins)	15	15	60
Time to repair node (t_repair) (hrs)	0.25	0.25	1
Time to init job (mins)	10	10	15
Network overhead (%)	0	0	0
Memory overhead (%)	0	5	0
<b>RESULTS</b>			
Cluster MTBF	4.8h	4.8h	2.9h
Interruptions/mo	149.3	149.3	248.8
Downtime loss	5.53%	5.53%	20.91%
Idle spare cost	0.62%	—	—
Performance overhead	0.00%	5.00%	0.00%
<b>Total Goodput Loss</b>	<b>6.14%</b>	<b>10.53%</b>	<b>20.91%</b>

# INFERENCEX

Open-source continuous inference benchmark — the bridge between raw hardware and \$/Mtok

## WHAT IT MEASURES

- Interactivity (tok/s/user) and per-GPU throughput (tok/s/GPU)
- P99 time-to-first-token (TTFT) across concurrency levels
- Tokens per megawatt (tok/s/MW) and joules per token
- \$ per million tokens at hyperscaler, neocloud, and rental price

*Every datapoint is a public GitHub Actions run — click any chart point to inspect the recipe, logs, and artifacts. Weekly DB releases for audit.*

## COVERAGE

<b>NVIDIA</b>	H100, H200, B200, B300, GB200, GB300
<b>AMD</b>	MI300X, MI325X, MI355X
<b>Models</b>	DeepSeek-R1, gpt-oss-120B, Llama-3.3-70B, Qwen, Kimi-K2, GLM-5, MiniMax-M2
<b>Frameworks</b>	vLLM, SGLang, TRT-LLM, Dynamo (TRT/vLLM/SGLang)
<b>Precisions</b>	FP4, FP8, BF16, INT4
<b>Modes</b>	Disaggregated prefill/decode, MTP, wide expert parallelism for MoE

# FROM BENCHMARK TO ECONOMICS

InferenceX gives you tok/s/GPU and joules/token. Add price and SLA and you have \$/Mtok — the unit the next section is built on.

1

## BENCHMARK

InferenceX measures tok/s/GPU, P99 TTFT, joules/token on real hardware across all models, frameworks, precisions.

2

## PRICE

Apply the GPU rental price (from ClusterMAX-rated providers) or hyperscaler list price to the throughput.

3

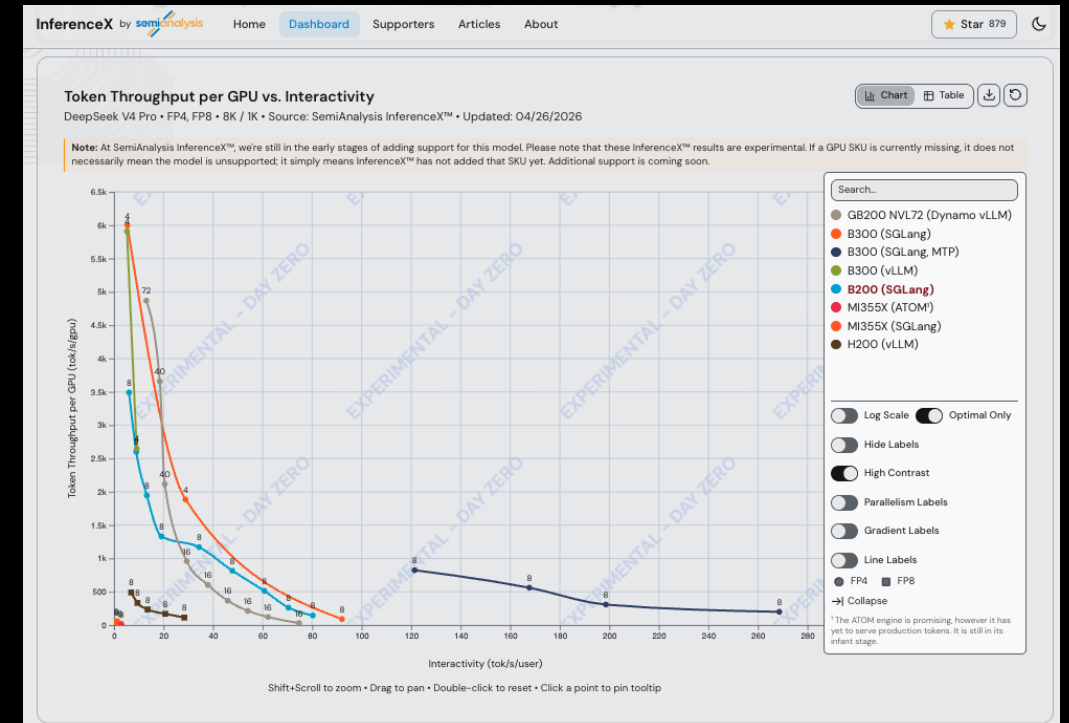
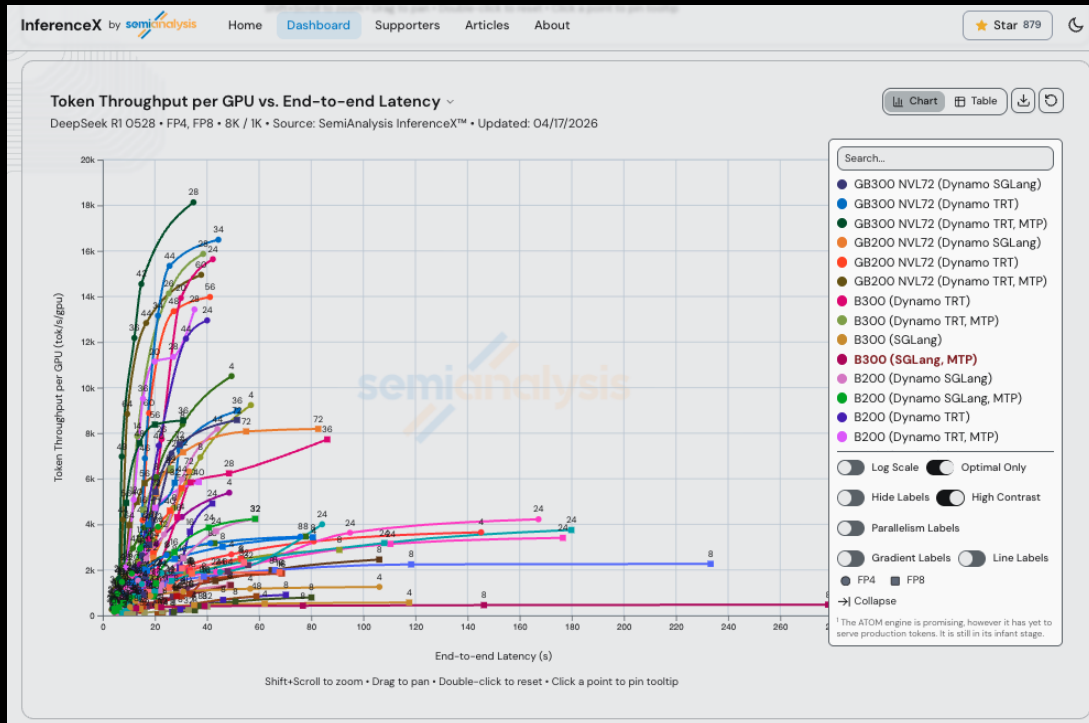
## \$/MTOK

Result: cost per million tokens at your SLA — the only unit the customer actually feels and pays for.

*Up next — tokenomics: who pays the \$/Mtok*

# INTRODUCTION TO INFERENCEX

An open-source, automated benchmark that moves at the same speed as the AI software ecosystem



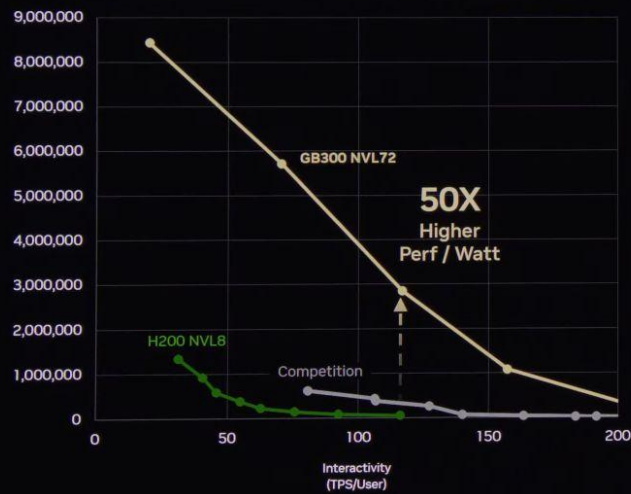
featuring day zero support for DeepSeek V4!

# INTRODUCTION TO INFERENCEX

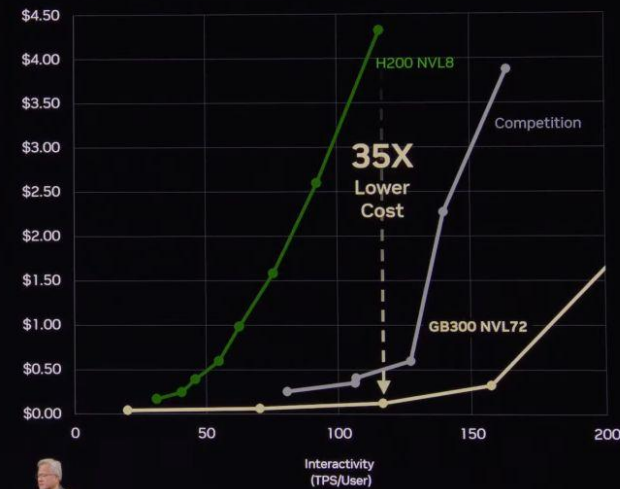
Featured on stage at the GTC 2026 Keynote by Inference King Jensen Huang

## NVIDIA Extreme Co-Design Revolutionized Token Cost “GB NVL72 Inference King”

Tokens per Watt Drives Factory Revenue



Performance Drives Token Cost



InferenceX  
by semianalysis



DeepSeek R1 0528 · FP4 · 1K/1K · Source: SemiAnalysis InferenceX



# UNDERSTANDING INFERENCEX

The screenshot shows the GitHub repository for InferenceX. The repository is public and has 1,164 commits, 149 forks, and 880 stars. The README is visible, titled "InferenceX™, Open Source Continuous Inference Standard and Research Platform". It mentions that it is trusted by operators of trillion-dollar token factories such as OpenAI, Microsoft, Oracle, etc., and ML community such as PyTorch Foundation, vLLM, SGLang, and Tri Dao. The repository is licensed under Apache 2.0. The repository structure includes folders like .claude, .github, benchmarks, experimental, runners, and utils, and files like .gitignore, .mcp.json, AGENTS.md, LICENSE, README.md, and perf-changelogyaml.

Software evolves every single day, delivering **continuous performance gains** on top of step jumps in hardware

AI software like SGLang, vLLM, TensorRT-LLM, CUDA and ROCm increase the **pareto frontier of performance** in incremental releases that are just **days apart**

Benchmarks conducted at a fixed point in time quickly go stale and are not representative of what can be achieved with the latest software packages

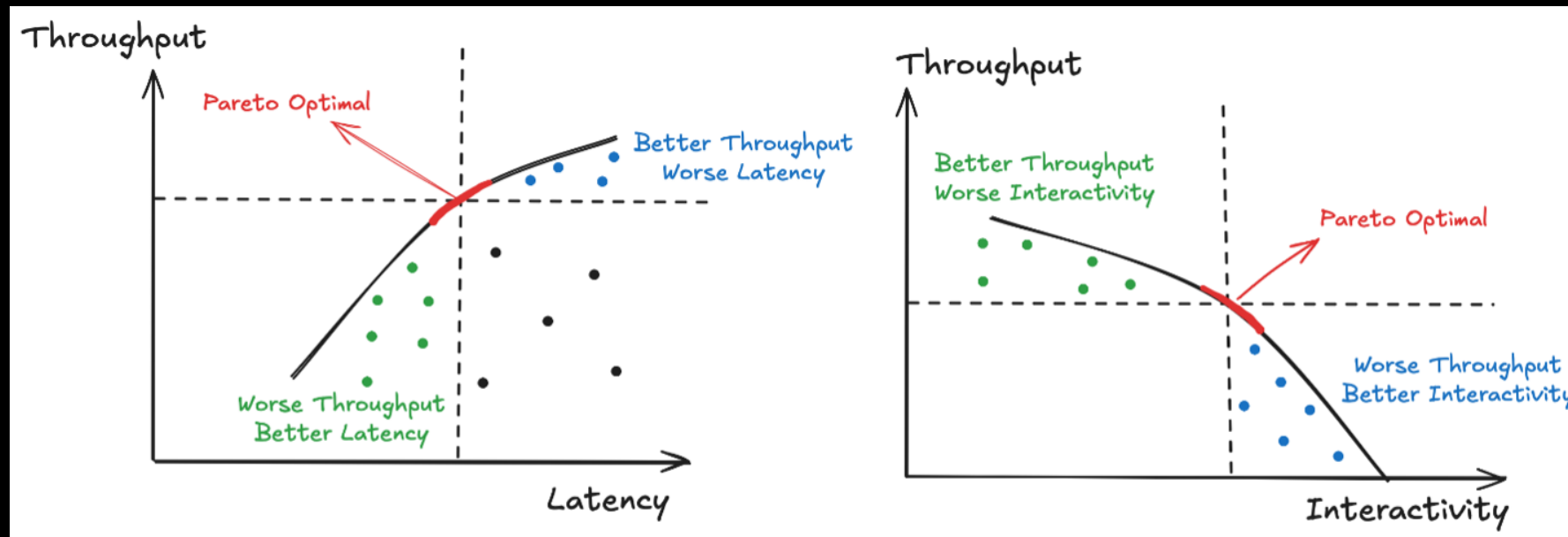
# UNDERSTANDING INFERENCECX: THE PARETO FRONTIER

## What It Is

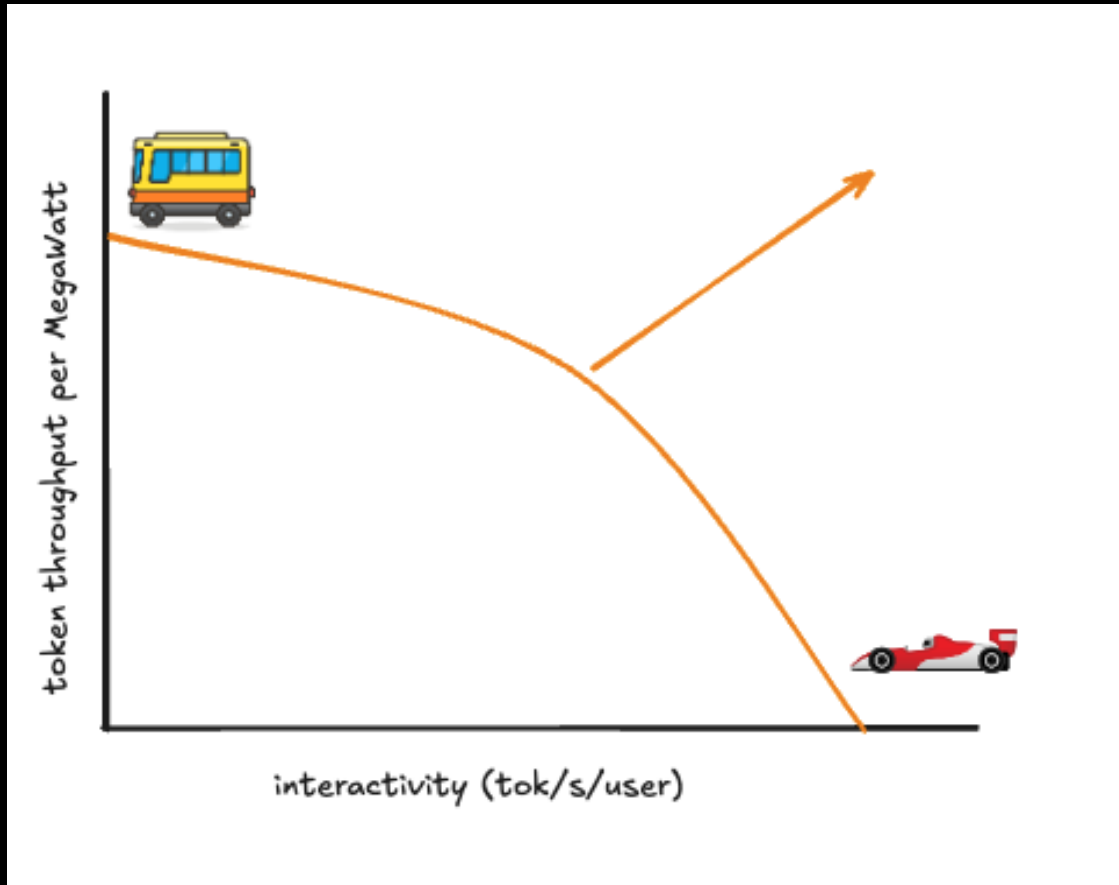
- Curve showing **best possible** throughput vs latency trade-offs
- **Pareto optimal point**: No other point improves one metric without sacrificing the other
- Points below frontier = **suboptimal**, wasting performance

## Why It Matters

- Operating below frontier = wasted resources = **lost revenue**
- Moving to frontier = **serve more users** OR **deliver faster responses** with same infrastructure
- Frontier = **maximum revenue potential**



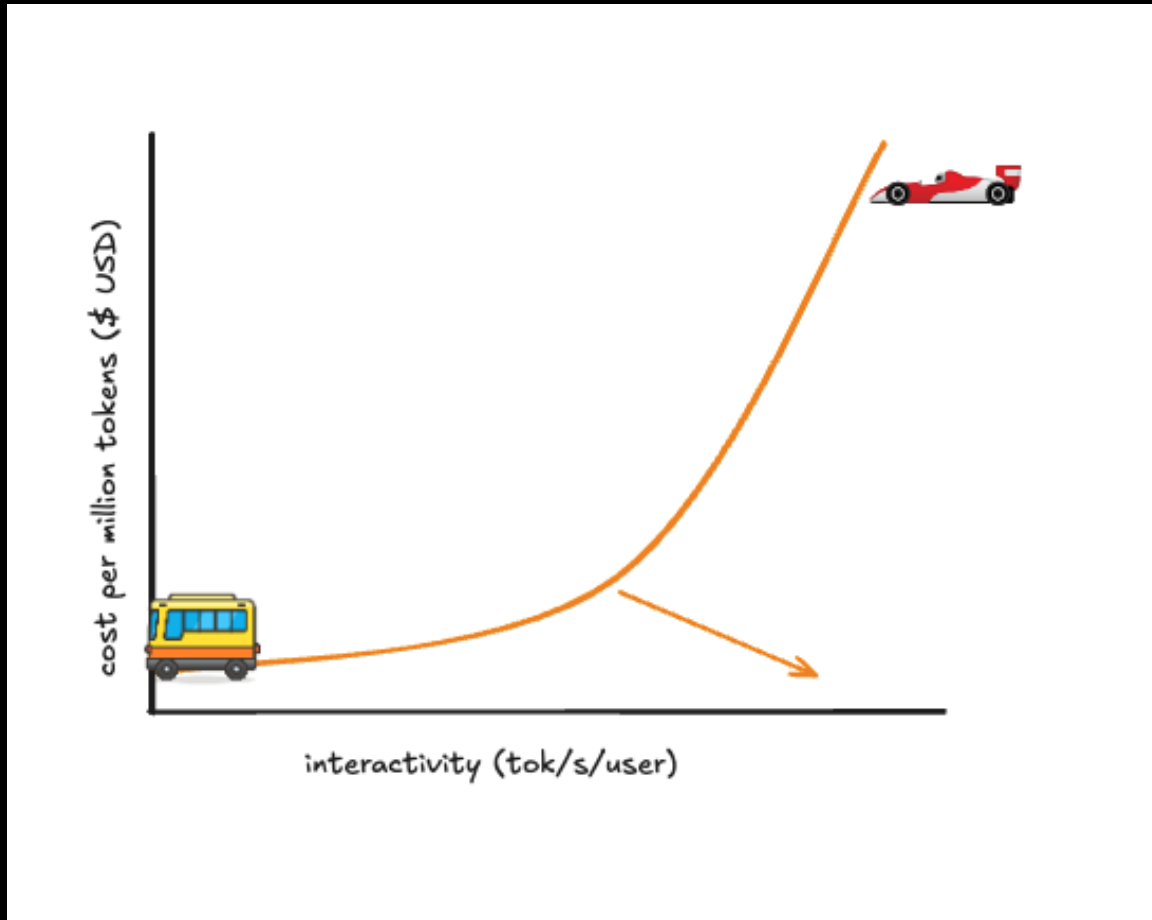
# TOKEN THROUGHPUT VS INTERACTIVITY (TOK/S/USER)



## Why It Matters

- Electricity & colocation cost make up to **20% of TCO**
- 20% lower tok/s/MW will only impact **4% of TCO**
- Higher interactivity → fewer tokens per datacenter
- Lower interactivity → higher tokens per datacenter

# TCO PER MILLION TOKENS VS INTERACTIVITY (TOK/S/USER)



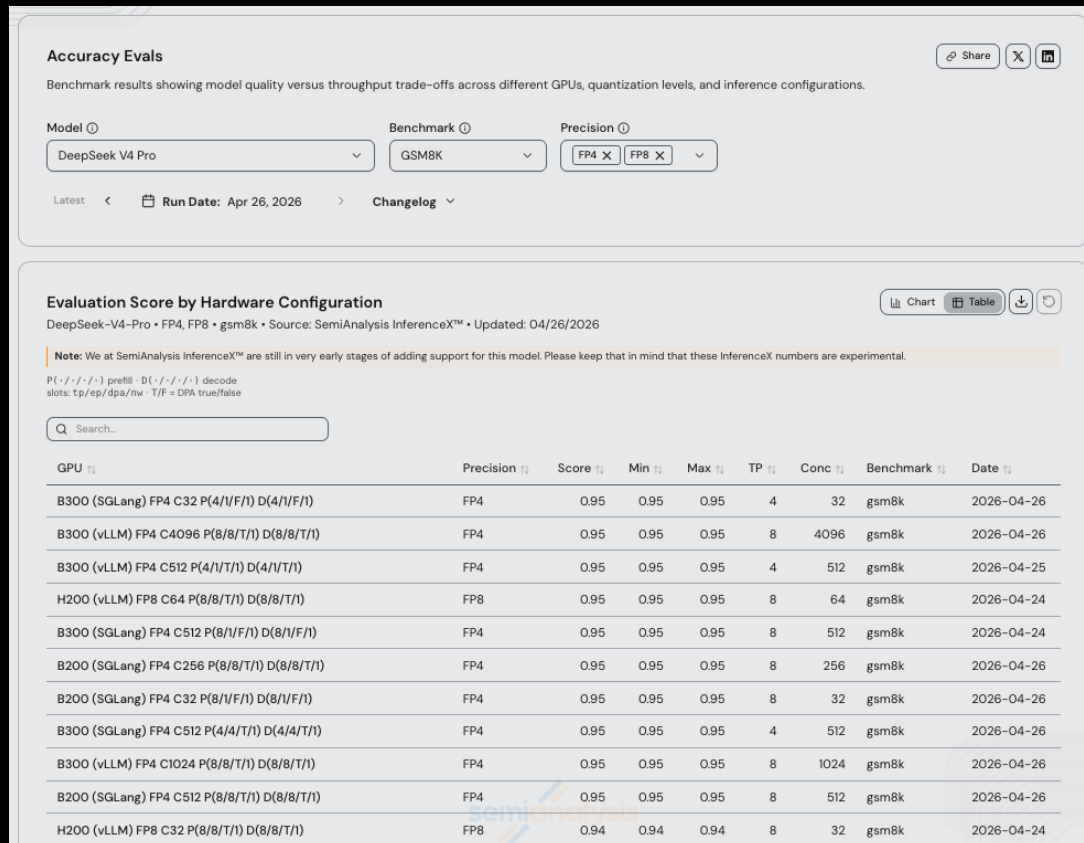
## Why It Matters

- Higher interactivity → fewer tokens/hour → higher \$/token
- Higher throughput → more tokens/hour → lower \$/token
- **Premium responsiveness** requires **premium pricing**

# ADDITIONAL FEATURES: GOING BEYOND PERFORMANCE

## Accuracy Evals

Confirming that performance improvements do not sacrifice model quality



## Historical Trends

At iso-interactivity, demonstrating throughput improvements over time



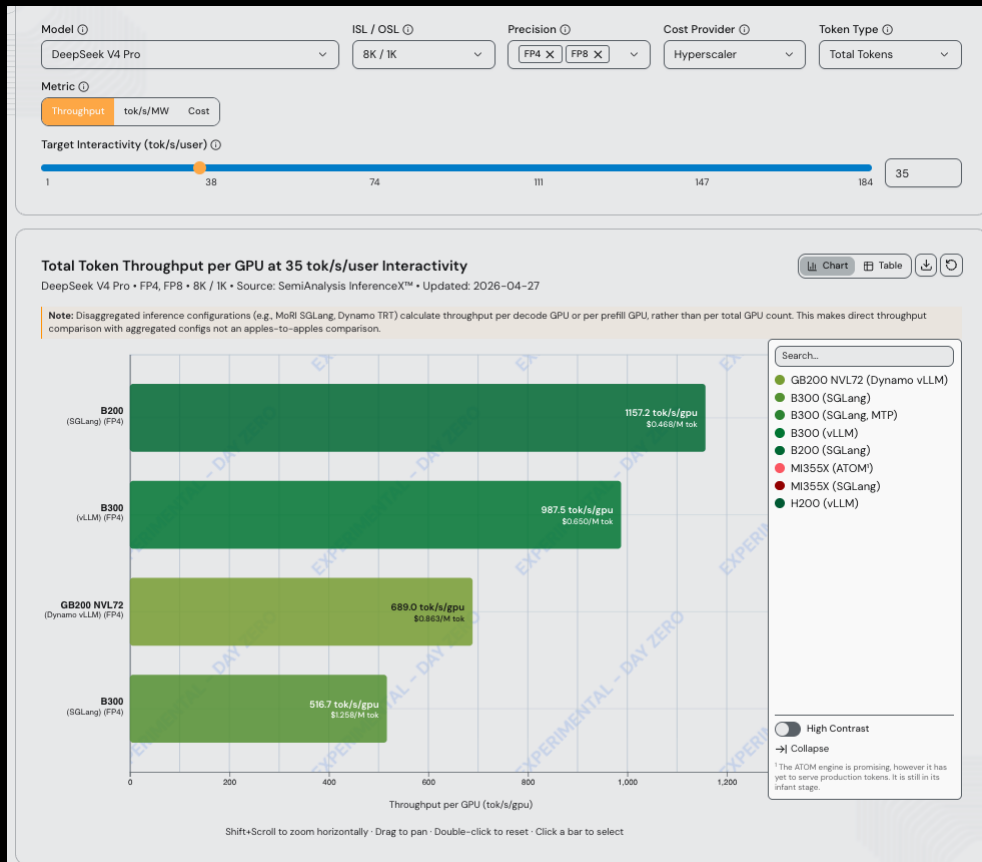
# ADDITIONAL FEATURES: GOING BEYOND PERFORMANCE

## TCO Calculator

Assessing which hardware is best for a given model

## Interactive Diagrams

Understanding GPU Systems, Model Architectures and More



# TOKENOMICS: WHERE DOES THE VALUE GO?

The AI Token Factory Economics Stack				
Layer	Units	H200	Business	SemiAnalysis Model
<b>Total Server Capital Cost</b>	<b>USD per 8 GPU</b>	<b>\$246,391</b>	Chip	Accelerator and AI TCO Model
Useful Life	Yrs	4.0	IAAS	AI TCO Model
<b>Depreciation Expense/GPU</b>	<b>USD/hr/GPU</b>	<b>\$0.88</b>	IAAS	
Colocation Cost	USD/kW/mth	\$110	IAAS	AI Datacenter Model
Electricity Cost	USD/kWh	\$0.087	IAAS	
<b>Operating Cost per GPU</b>	<b>USD/hr/GPU</b>	<b>\$0.35</b>	IAAS	
<b>Total Variable Cost per GPU</b>	<b>USD/hr/GPU</b>	<b>\$1.23</b>	IAAS	AI TCO Model
<b>Revenue per GPU</b>	<b>USD/hr/GPU</b>	<b>\$1.90</b>	IAAS	
<i>Neocloud Gross Margin</i>	%	35%	IAAS	
Effective Inference Throughput/GPU <sup>1</sup>	Tokens/s/GPU	735	PAAS	InferenceMAX™
<b>Inference Compute Cost</b>	<b>USD/M Tok</b>	<b>\$0.72</b>	PAAS	
<b>Inference API Blended Price</b>	<b>USD/M Tok</b>	<b>\$1.75</b>	PAAS	
<i>Model Gross Margin</i>	%	59.0%	PAAS	
Tokens/User/Month	Tok/user/mth	10,000,000	Application	AI Tokenomics Model
3P Token Unit Cost	USD/M Tok	\$1.75	Application	
Average Revenue per User	USD/mth	\$20.00	Application	
<b>3P Gross profit per User/Mth</b>	<b>USD/mth</b>	<b>\$2.50</b>	Application	
<i>3P Application Gross Margin</i>	%	12.5%	Application	
1P Token Unit Cost	USD/M Tok	\$0.72		AI Tokenomics Model
<b>1P Gross profit per User/Mth</b>	<b>USD/mth</b>	<b>\$12.82</b>	Application	
<i>1P Application Gross Margin</i>	%	64.1%	Application	

1. DeepSeek R1 FP8 on H200 using SGLang. Uses 8k input, 1k output tokens, 43 interactivity.

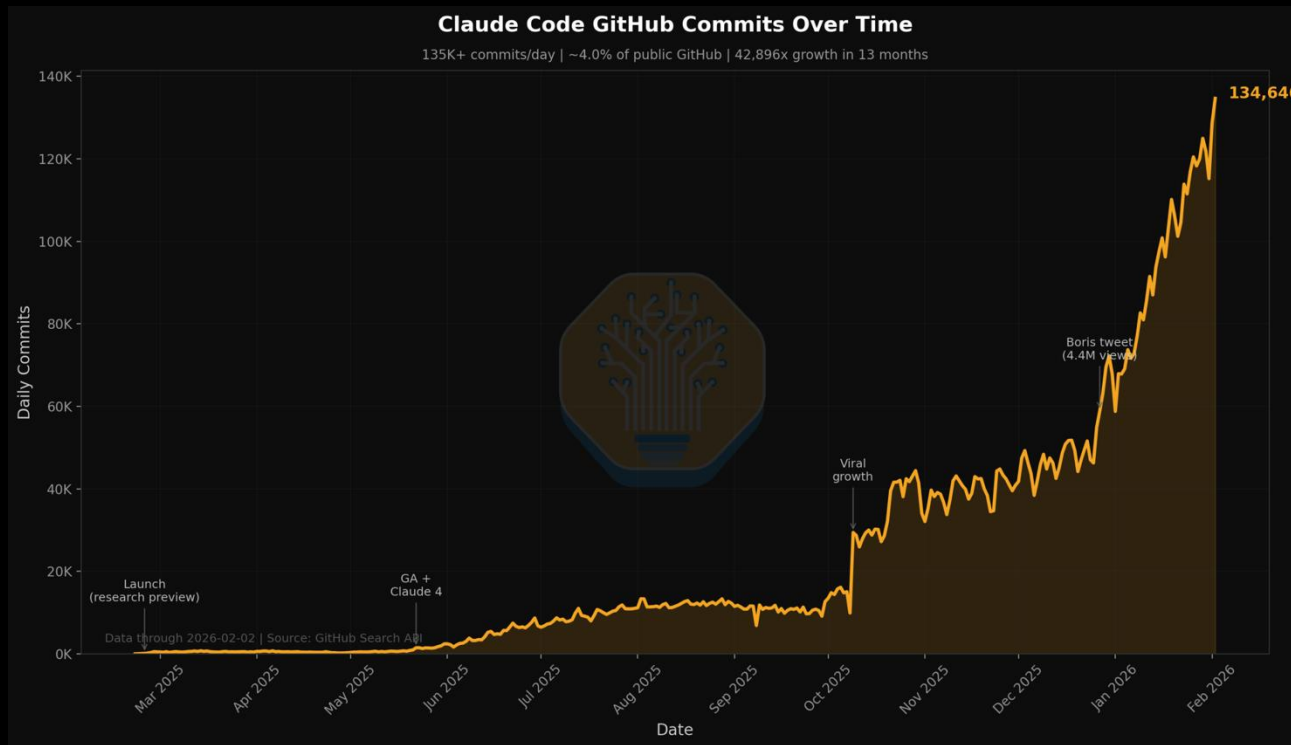
PART 04

# End User Demand

---

Tokenomics: who actually uses the tokens, how it is monetized, and whether the revenue justifies the buildout.

# CLAUDE CODE HAS ARRIVED



## Claude Code is the Inflection Point

What It Is, How We Use It, Industry Repercussions, Microsoft's Dilemma, Why Anthropic Is Winning

DOUG O'LAUGHLIN, JEREMIE ELIAHOU ONTIVEROS, JORDAN NANOS, AND 2 OTHERS  
FEB 05, 2026 · PAID

538 6 81

Share

4% of GitHub public commits are being authored by Claude Code right now. At the current trajectory, we believe that Claude Code will be 20%+ of all daily commits by the end of 2026. While you blinked, AI consumed all of software development.

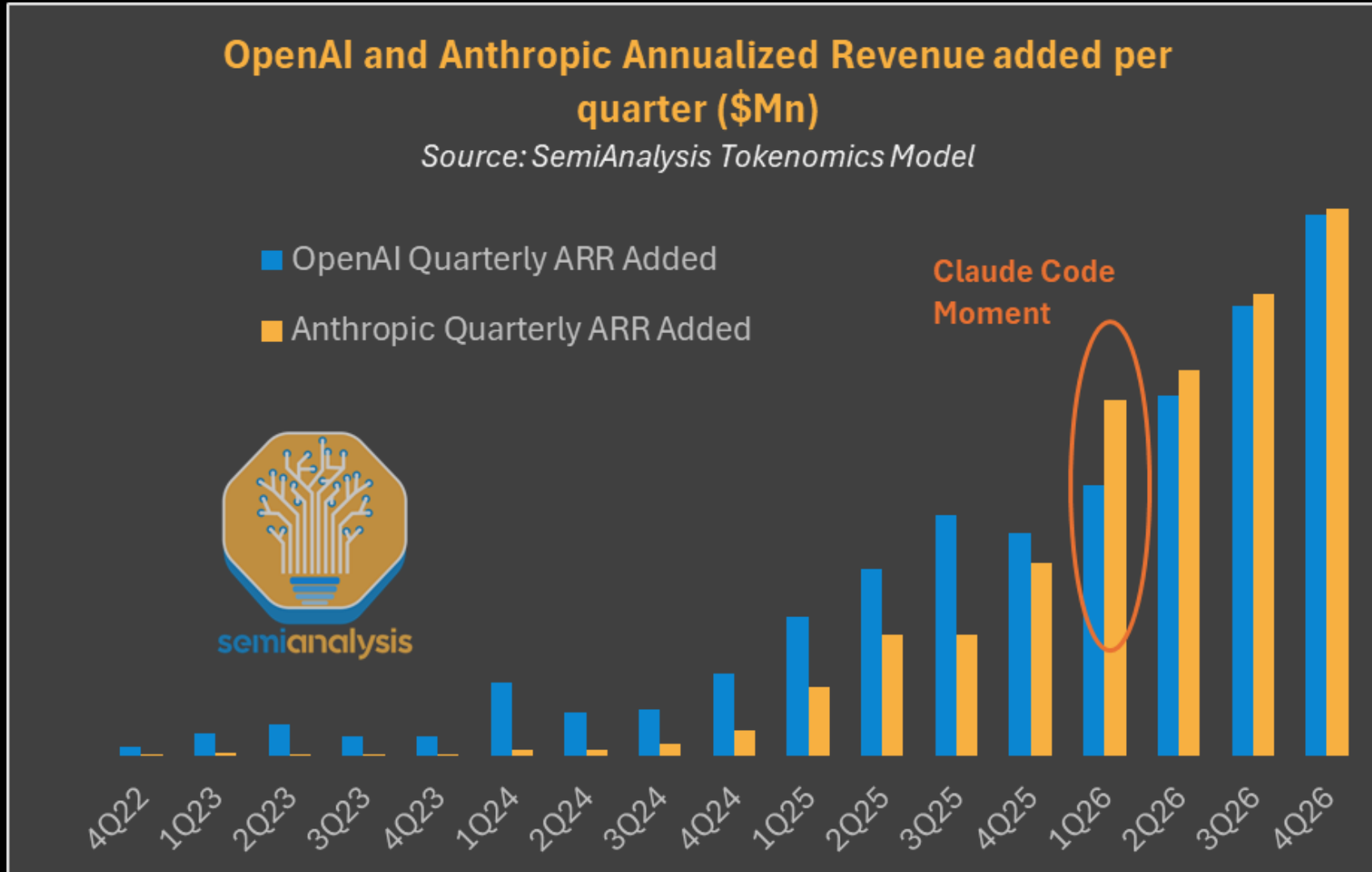
Our sister publication Fabricated Knowledge described software like linear TV during the rise of the internet and thinks that the rise of Claude Code is going to be a new layer of intelligence on top of software akin to DRAM versus NAND. Today SemiAnalysis is going to dive into the repercussions of Claude Code, what it is, and why Claude is so good.

Source: <https://newsletter.semianalysis.com/p/claude-code-is-the-inflection-point>

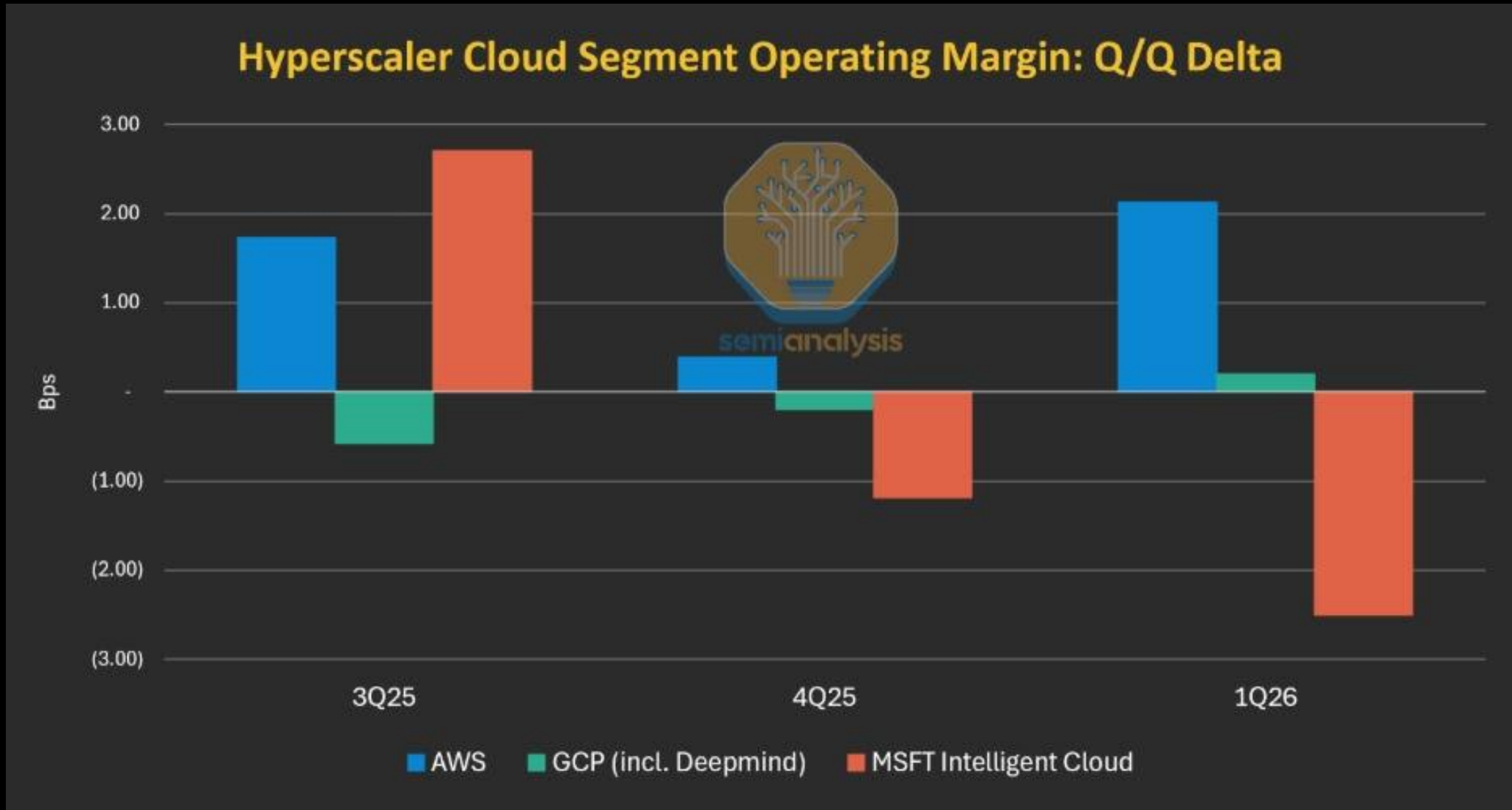
# ANTHROPIC GROWTH AS A RESULT

## OpenAI and Anthropic Annualized Revenue added per quarter (\$Mn)

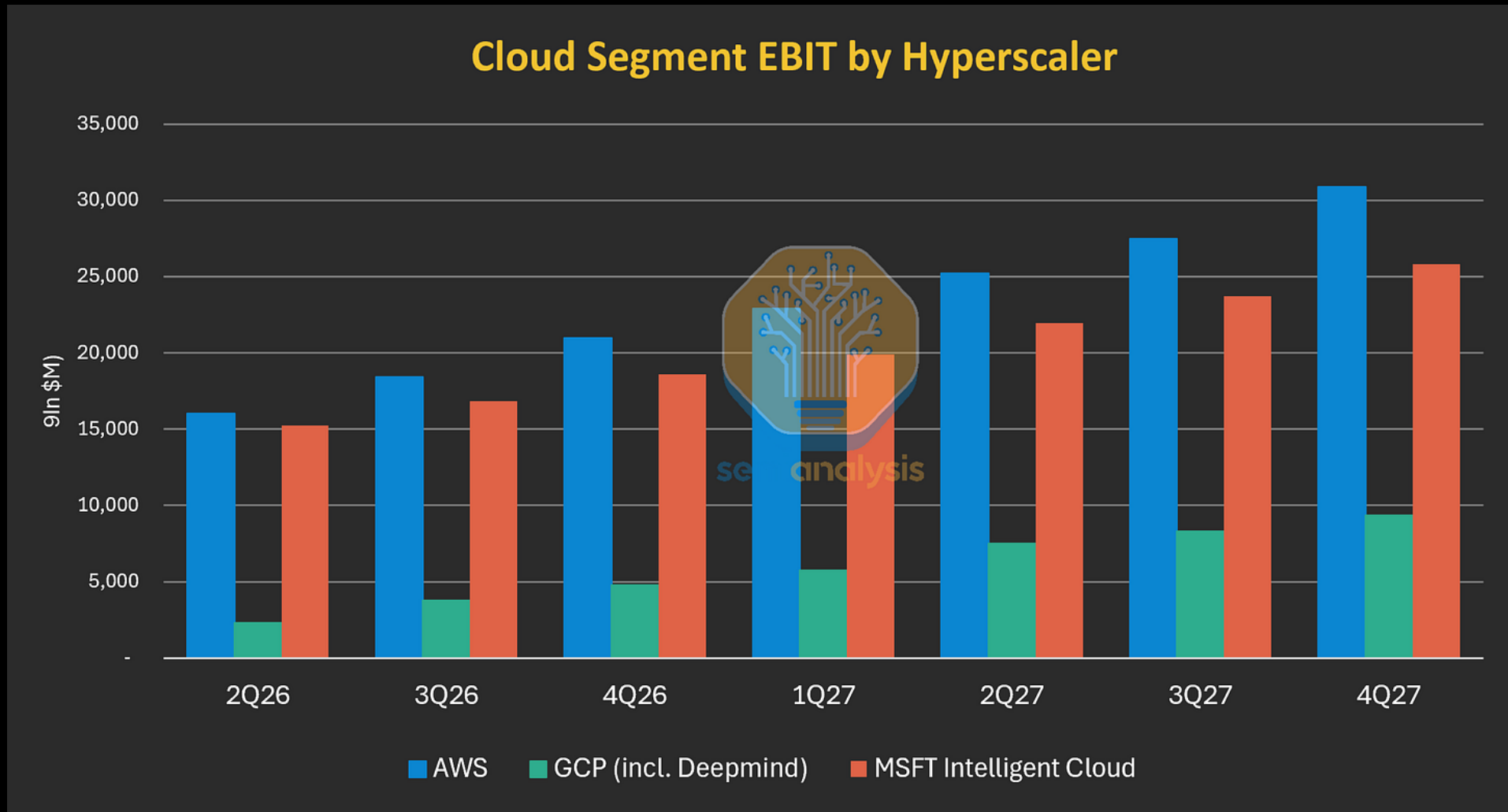
Source: SemiAnalysis Tokenomics Model



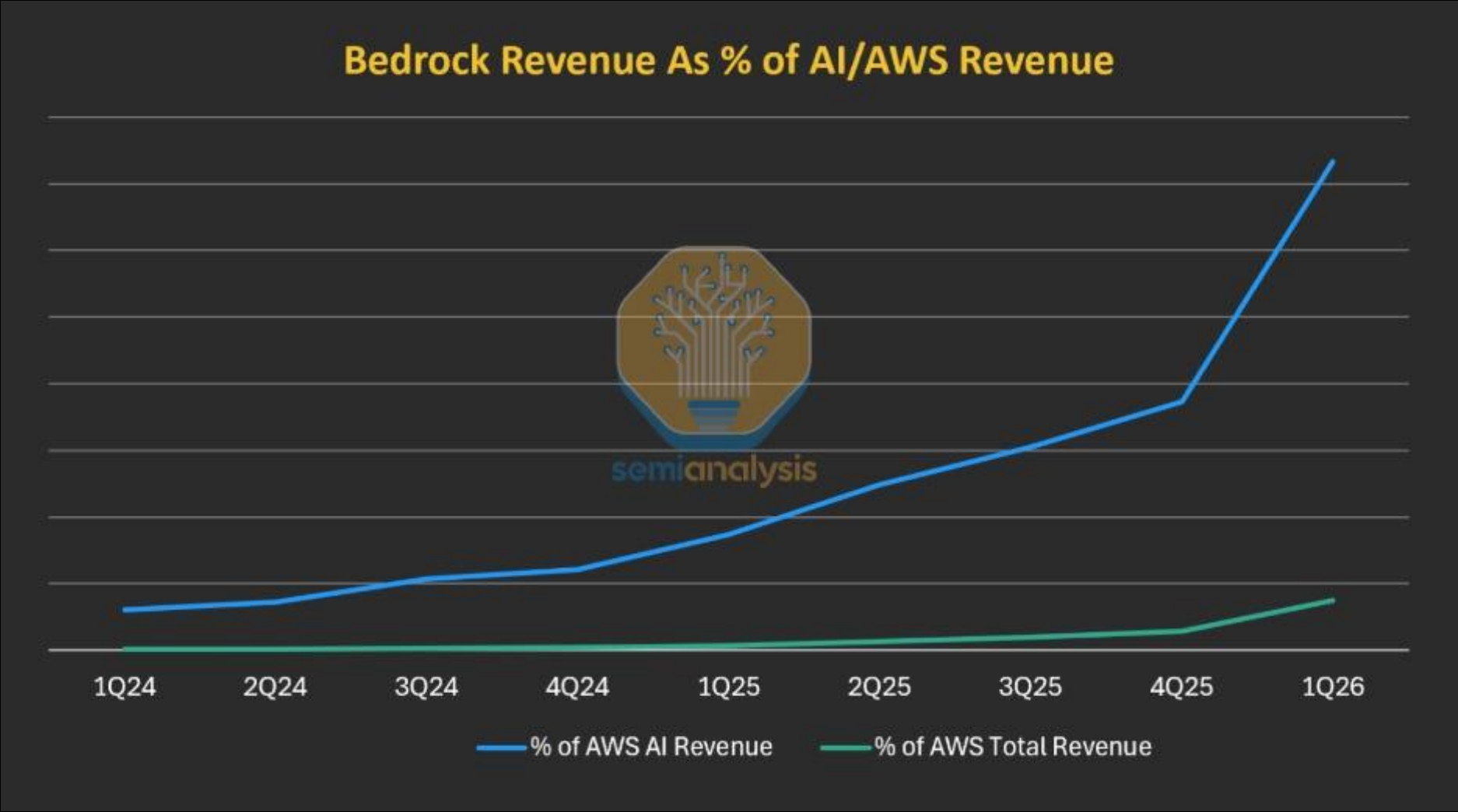
# AMAZON IS A KEY BENEFICIARY OF THIS GROWTH



# AMAZON IS A KEY BENEFICIARY OF THIS GROWTH

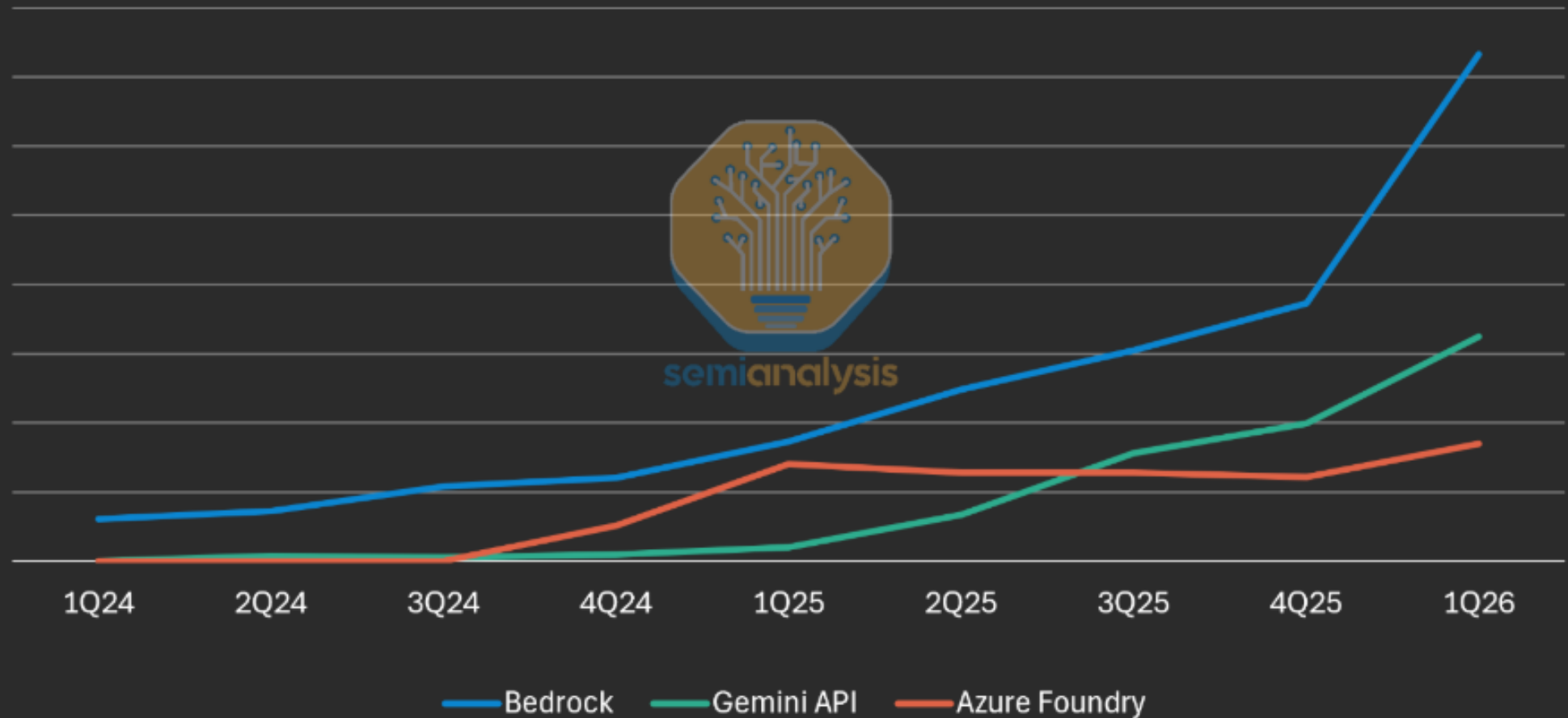


# BEDROCK IS GROWING LIKE CRAZY

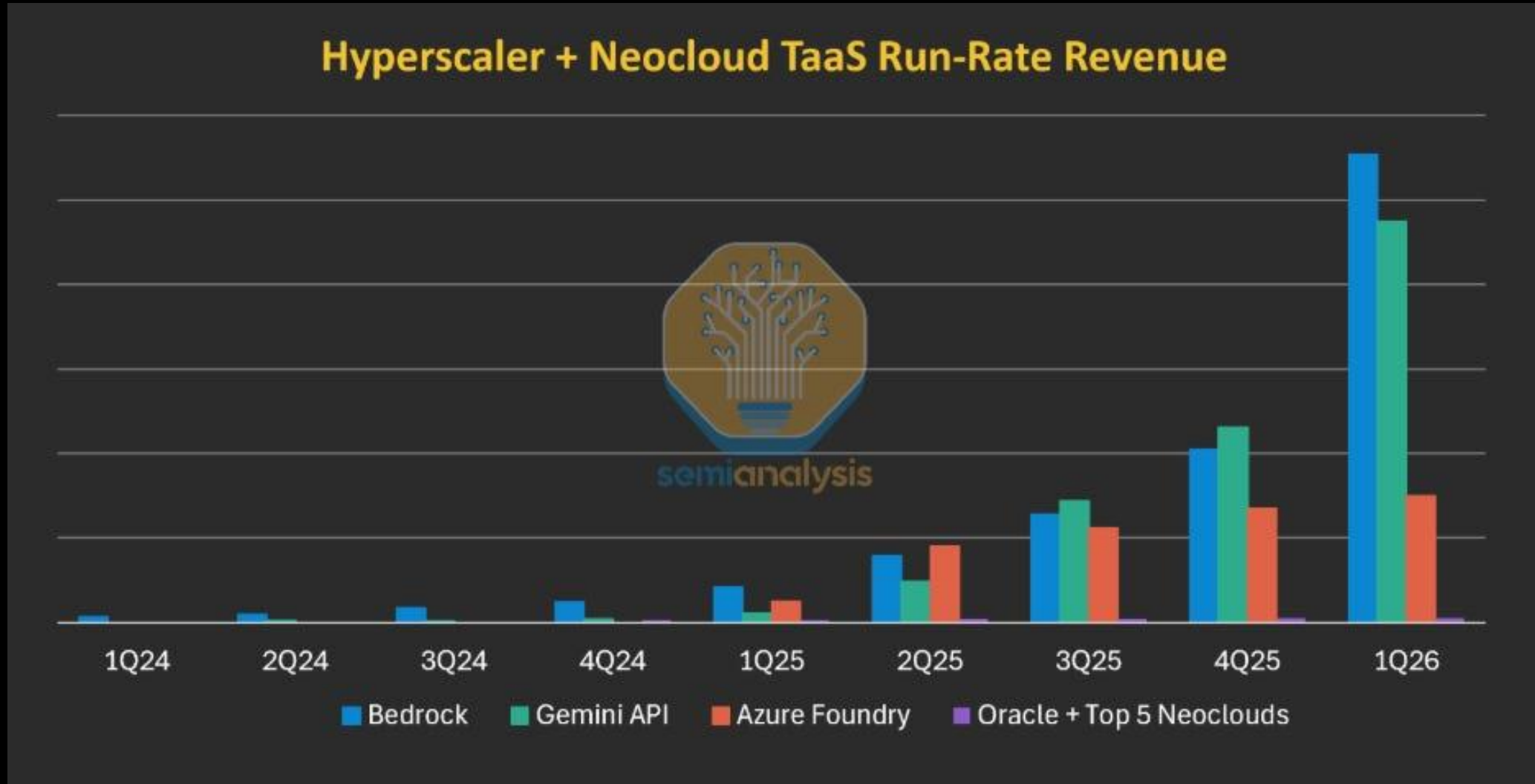


# TOKEN GROWTH IS STRONG

## Token-as-a-Service Revenue as % of AI Revenue



# NEOCLOUDS ARE LARGELY IRRELEVANT



# WE EXPECT THIS TO CONTINUE



# LOOK AT HOW FAST THEY'RE BUILDING



AWS / Anthropic AI Training Datacenters as of Q3 2025

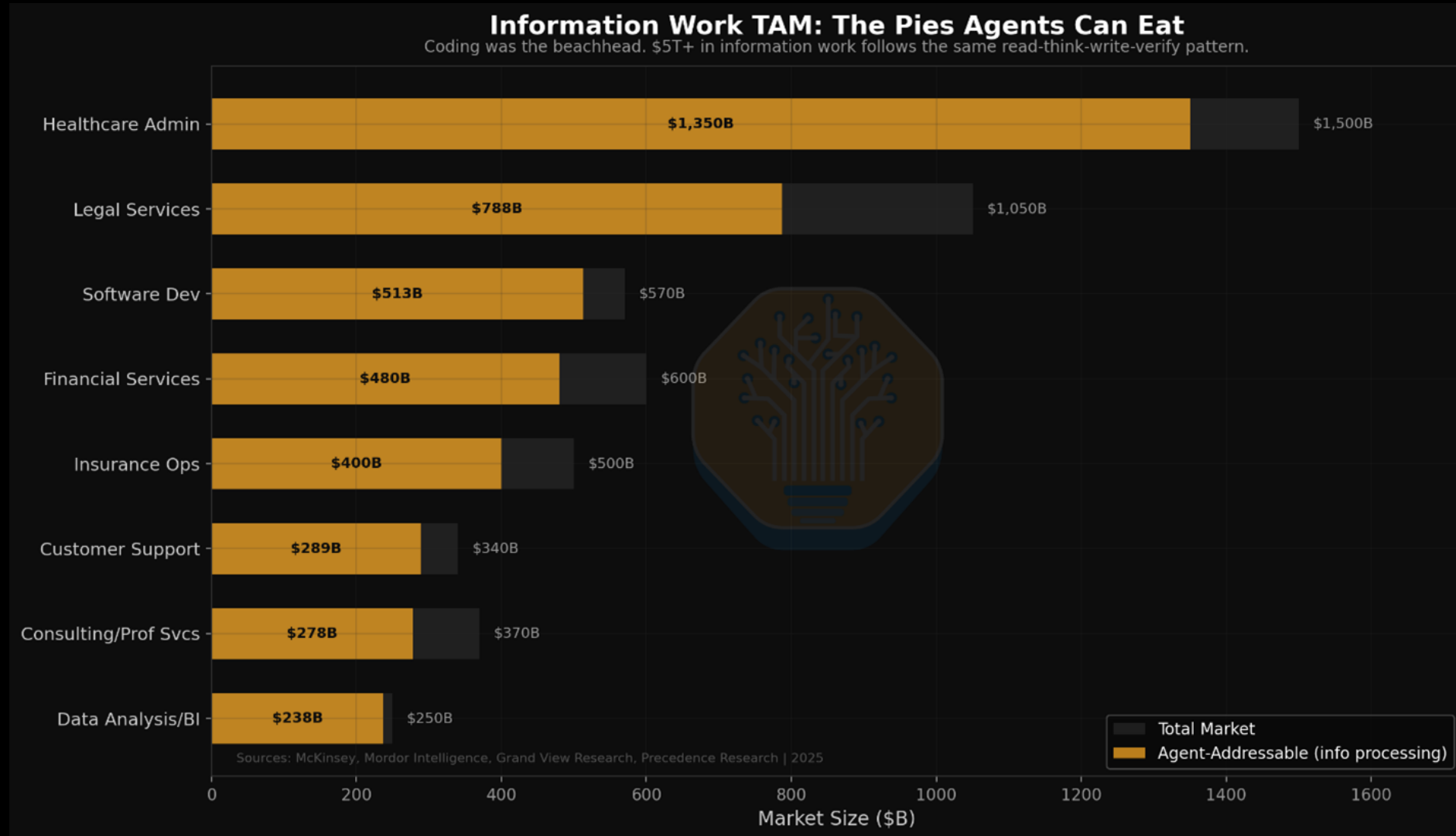
Source: SemiAnalysis Datacenter Industry Model



AWS / Anthropic AI Training Datacenters as of Q3 2024

Source: SemiAnalysis Datacenter Industry Model

# AGENTS ARE JUST GETTING STARTED



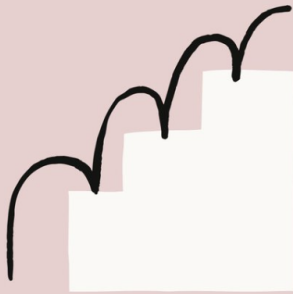
# OPENAI & ANTHROPIC ARE HIGHLY PROFITABLE

**Anthropic** @AnthropicAI

Earlier this month, our run-rate revenue crossed \$47 billion.

This growth has been driven by organizations across many industries deploying Claude in their core operations, and by a growing number of people using it for their everyday work.

Read more:



Anthropic raises \$65B in Series H funding at \$965B post-money valuation

From anthropic.com

source: <https://www.anthropic.com/news/series-h>

OpenAI was the fastest technology platform to reach 10 million users, the fastest to 100 million users, and soon the fastest to 1 billion weekly active users. Within a year of launching ChatGPT, we reached \$1B in revenue. By the end of 2024 we were generating \$1B per quarter. We are now generating **\$2B in revenue per month**. At this stage, we are growing revenue four times faster than the companies who defined the Internet and mobile eras, including Alphabet and Meta.

source: <https://openai.com/index/accelerating-the-next-phase-ai/>

← **Post**

**Joey Brookhart** @SaasquatchC

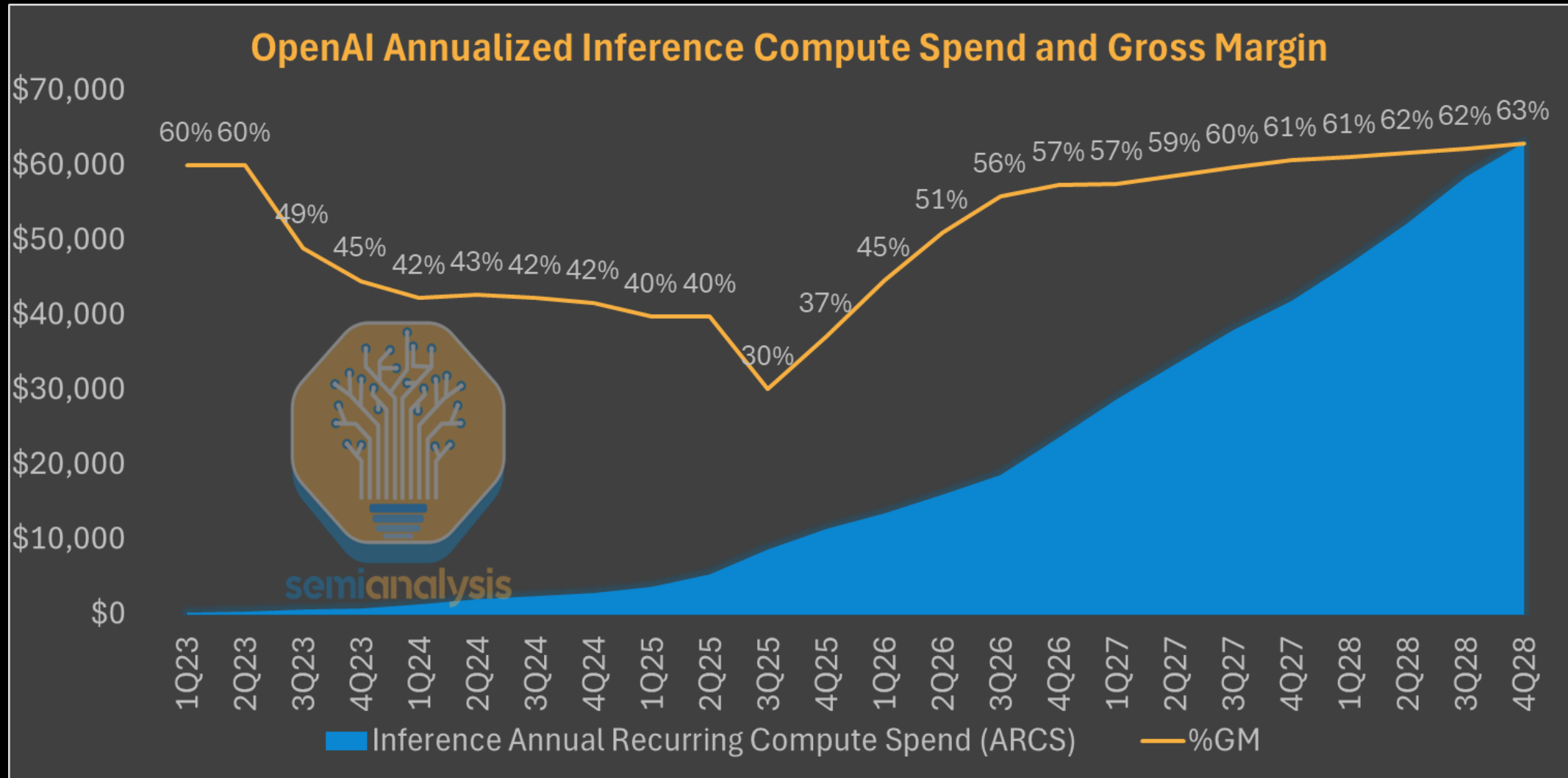
Congrats to our Tokenomics team at SemiAnalysis who had the Anthropic revenue numbers, cost buckets, and adjusted EBIT numbers dead on for Q1 and Q2

7:40 PM · May 20, 2026 · 20.3K Views

2 6 205 47

Relevant ▾ View quotes >

# THE “SHIFT FROM TRAINING TO INFERENCE” IS A LIE



# THANK YOU

QUESTIONS?

---

## JORDAN NANOS

*Member of Technical Staff · SemiAnalysis*

**Newsletter**      [newsletter.semianalysis.com](https://newsletter.semianalysis.com)

**SemiAnalysis**      [semianalysis.com](https://semianalysis.com)

**ClusterMAX**      [clustermax.ai](https://clustermax.ai)

**InferenceX**      [inferencex.com](https://inferencex.com)

**Tokenomics**      [tokenomics.info](https://tokenomics.info)

## THE TAKEAWAY

- Supply chain (CoWoS, HBM) decides who gets chips.
- Power + cooling innovations defines where they land.
- Networking decides how they perform.
- Tokenomics decides who earns the ROIC.